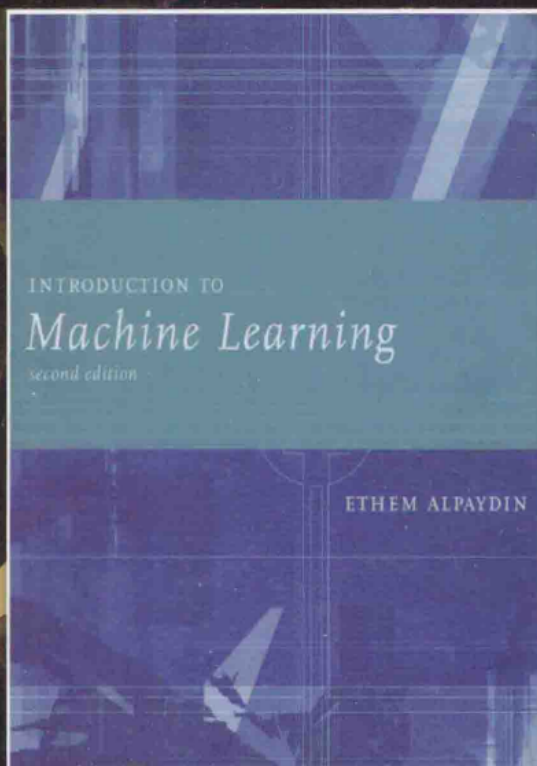


机器学习导论

(土耳其) Ethem Alpaydin 著 范明 咎红英 牛常勇 译

Introduction to Machine Learning
Second Edition



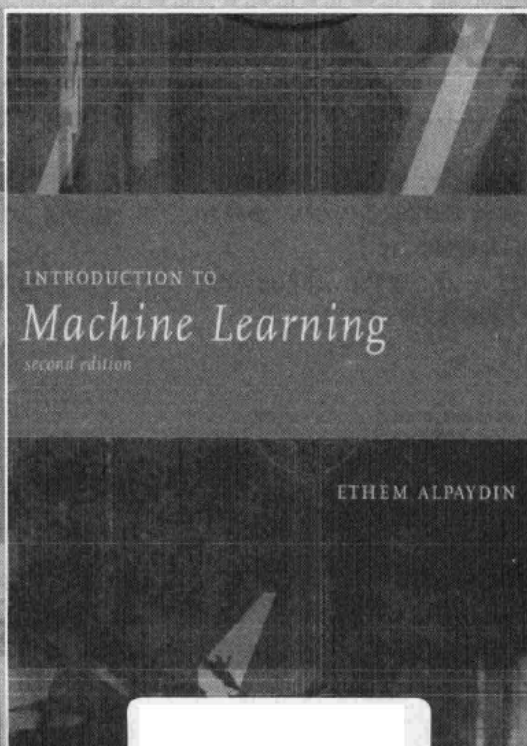
计 算 机 科 学 丛

原书第2版

机器学习导论

(土耳其) Ethem Alpaydin 著 范明 咎红英 牛常勇 译

Introduction to Machine Learning
Second Edition



机械工业出版社
China Machine Press

www.aibbt.com 让未来触手可及

图书在版编目 (CIP) 数据

机器学习导论 (原书第 2 版) / (土耳其) 阿培丁 (Alpaydin, E.) 著; 范明等译. —北京: 机械工业出版社, 2014.3

(计算机科学丛书)

书名原文: Introduction to Machine Learning, Second Edition

ISBN 978-7-111-45377-2

I. 机… II. ①阿… ②范… III. 机器学习—研究 IV. TP181

中国版本图书馆 CIP 数据核字 (2013) 第 321864 号

本书版权登记号: 图字: 01-2013-1430

Ethem Alpaydin: Introduction to Machine Learning, Second Edition (ISBN 978-0-262-01243-0).

Original English language edition copyright © 2010 by Massachusetts Institute of Technology.

Simplified Chinese Translation Copyright © 2014 by China Machine Press.

Simplified Chinese translation rights arranged with MIT Press through Bardon-Chinese Media Agency.

No part of this book may be reproduced or transmitted in any form or by any means, electronic or mechanical, including photocopying, recording or any information storage and retrieval system, without permission, in writing, from the publisher.

All rights reserved..

本书中文简体字版由 MIT Press 通过 Bardon-Chinese Media Agency 授权机械工业出版社在中华人民共和国境内独家出版发行。未经出版者书面许可, 不得以任何方式抄袭、复制或节录本书中的任何部分。

本书讨论了机器学习在统计学、模式识别、神经网络、人工智能、信号处理等不同领域的应用, 其中涵盖了监督学习、贝叶斯决策理论、参数方法、多元方法、多层感知器、局部模型、隐马尔可夫模型、分类算法评估和比较以及增强学习。

本书可供完成计算机程序设计、概率论、微积分和线性代数课程的高年级本科生和研究生使用, 也可供对机器学习感兴趣的工程技术人员参考。

出版发行: 机械工业出版社 (北京市西城区百万庄大街 22 号 邮政编码: 100037)

责任编辑: 肖晓慧

印 刷: 蕺城市京瑞印刷有限公司

版 次: 2014 年 4 月第 1 版第 1 次印刷

开 本: 185mm × 260mm 1/16

印 张: 22

书 号: ISBN 978-7-111-45377-2

定 价: 59.00 元

凡购本书, 如有缺页、倒页、脱页, 由本社发行部调换

客服热线: (010) 88378991 88361066

投稿热线: (010) 88379604

购书热线: (010) 68326294 88379649 68995259

读者信箱: hzjsj@hzbook.com

版权所有·侵权必究

封底无防伪标均为盗版

本书法律顾问: 北京大成律师事务所 韩光 / 邹晓东

文艺复兴以降，源远流长的科学精神和逐步形成的学术规范，使西方国家在自然科学的各个领域取得了垄断性的优势；也正是这样的传统，使美国在信息技术发展的六十多年间名家辈出、独领风骚。在商业化的进程中，美国的产业界与教育界越来越紧密地结合，计算机学科中的许多泰山北斗同时身处科研和教学的最前线，由此而产生的经典科学著作，不仅擘划了研究的范畴，还揭示了学术的源变，既遵循学术规范，又自有学者个性，其价值并不会因年月的流逝而减退。

近年，在全球信息化大潮的推动下，我国的计算机产业发展迅猛，对专业人才的需求日益迫切。这对计算机教育界和出版界既是机遇，也是挑战；而专业教材的建设在教育战略上显得举足轻重。在我国信息技术发展时间较短的现状下，美国等发达国家在其计算机科学发展的几十年间积淀和发展的经典教材仍有许多值得借鉴之处。因此，引进一批国外优秀计算机教材将对我国计算机教育事业的发展起到积极的推动作用，也是与世界接轨、建设真正的世界一流大学的必由之路。

机械工业出版社华章公司较早意识到“出版要为教育服务”。自1998年开始，我们就将工作重点放在了遴选、移译国外优秀教材上。经过多年的不懈努力，我们与 Pearson, McGraw-Hill, Elsevier, MIT, John Wiley & Sons, Cengage 等世界著名出版公司建立了良好的合作关系，从他们现有的数百种教材中甄选出 Andrew S. Tanenbaum, Bjarne Stroustrup, Brian W. Kernighan, Dennis Ritchie, Jim Gray, Alfred V. Aho, John E. Hopcroft, Jeffrey D. Ullman, Abraham Silberschatz, William Stallings, Donald E. Knuth, John L. Hennessy, Larry L. Peterson 等大师名家的一批经典作品，以“计算机科学丛书”为总称出版，供读者学习、研究及珍藏。大理石纹理的封面，也正体现了这套丛书的品位和格调。

“计算机科学丛书”的出版工作得到了国内外学者的鼎力相助，国内的专家不仅提供了中肯的选题指导，还不辞劳苦地担任了翻译和审校的工作；而原书的作者也相当关注其作品在中国的传播，有的还专程为其书的中译本作序。迄今，“计算机科学丛书”已经出版了近两百个品种，这些书籍在读者中树立了良好的口碑，并被许多高校采用为正式教材和参考书籍。其影印版“经典原版书库”作为姊妹篇也被越来越多实施双语教学的学校所采用。

权威的作者、经典的教材、一流的译者、严格的审校、精细的编辑，这些因素使我们的图书有了质量的保证。随着计算机科学与技术专业学科建设的不断完善和教材改革的逐渐深化，教育界对国外计算机教材的需求和应用都将步入一个新的阶段，我们的目标是尽善尽美，而反馈的意见正是我们达到这一终极目标的重要帮助。华章公司欢迎老师和读者对我们的工作提出建议或给予指正，我们的联系方式如下：

华章网站：www.hzbook.com

电子邮件：hzjsj@hzbook.com

联系电话：(010) 88379604

联系地址：北京市西城区百万庄南街1号

邮政编码：100037



华章教育

华章科技图书出版中心

中文版序

Introduction to Machine Learning, Second Edition

机器学习是计算机科学发展最快的领域之一。每年，我们都看到新的应用，并且学习算法的理论也发展很快。

我十分高兴地看到我的书的第2版用中文出版，并就为翻译所做出的努力，感谢范明教授。在此之前，范教授翻译了本书第1版以及一些统计学习和数据挖掘的名著。

我希望我的书的中文读者发现本书是有益的，并且就像我乐于写它一样乐于阅读它。

Ethem Alpaydin

于伊斯坦布尔博阿齐奇大学

2013年12月

Preface of the Chinese Edition

Machine learning is one of the fastest developing fields in computer science. Every year, we are seeing new applications, and the theory of learning algorithms is also developing very fast.

It gives me great pleasure to see the second edition of my book printed in Chinese, and for the effort in doing the translation, I would like to thank Professor Fan who previously have translated the first edition as well as several other well-known texts on statistics and data mining.

I hope that the Chinese readers of my book will find it beneficial and enjoy reading it as much as I enjoyed writing it.

Ethem Alpaydin

Boğazici University, Istanbul

December 2013

自从有计算机以来，人们就希望计算机能够学习。然而，机器学习真正取得实质性进展，能够成功地解决一些实际问题，并最终成为一个学科分支还是近 20 余年的事。

对于许多问题，我们的前人和先行者已经知道如何求解。例如，欧几里得告诉我们可以用辗转相除法求两个整数的最大公约数；Dijkstra 告诉我们如何有效地求两点之间的最短路径；Hoare 向我们展示了怎样将杂乱无章的对象快速排序……对于这些问题，我们清楚地知道求解步骤。因此，让计算机求解这些问题只需要设计算法和数据结构、进行编程，而不需要让计算机学习。

还有一些问题，人们可以轻而易举地做好，但是却无法解释清楚我们是如何做的。例如，尽管桌子千差万别、用途各异，但是我们一眼就能看出某个物体是否是桌子；尽管不同的人的手写阿拉伯数字大小不一、笔画粗细不同，但是我们还是可以轻易识别一个数字是不是 8；尽管声音时大时小，有时可能还有点沙哑，但是我们还是可以不用费力气地听出熟人的声音。诸如此类的例子不胜枚举。对于这些问题，我们不知道求解步骤。因此，让计算机来做这些事就需要让计算机学习。

我们知道桌子不是木材和各种材料的随机堆砌，手写数字不是像素的随机分布，熟人的声音也不是各种声波的随机混合。现实世界总是有规律的。机器学习正是从已知实例中自动发现规律，建立对未知实例的预测模型；根据经验不断提高，不断改进预测性能。

这是一本全面论述机器学习这一主题的教科书，适合作为高等院校计算机相关专业高年级本科生和研究生机器学习入门课程的教材。该书涵盖了监督学习、贝叶斯决策理论、参数方法、多元方法、维度归约、聚类、非参数方法、决策树、线性判别式、多层感知器、局部模型、隐马尔可夫模型、分类算法评估和比较、组合多学习器以及增强学习。作者对来自统计学、模式识别、神经网络、人工智能、信号处理、控制和数据挖掘等不同领域的机器学习问题和学习方法进行了统一论述。

第 2 版从 16 章扩展到 19 章，除增加 3 章外，对许多章节都进行了改写和扩充，以便反映机器学习的新进展。尤其是，核方法、贝叶斯估计和图模型这三个在第 1 版只用几节介绍的主题都各自扩充成了一整章，进行了更深入的讨论。此外，第 2 版还更加全面地介绍了机器学习实验的设计与分析，这在同类书籍中是独具特色的。

第 2 版由范明翻译。咎红英（第 1~6 章）和牛常勇（第 15、17 和 18 章）参加了第 1 版的翻译。原书作者 Ethem Alpaydin 为第 2 版的中文版重新写了序。

译文中的错误和不当之处，敬请读者朋友指正。意见和建议请发往 mfan@zzu.edu.cn，我们不胜感激。

范 明

2013 冬于郑州大学

机器学习使用实例数据或过去的经验训练计算机，以优化性能标准。当人们不能直接编写计算机程序解决给定的问题，而是需要借助于实例数据或经验时，就需要学习。一种需要学习的情况是人们没有专门技术，或者不能解释他们的专门技术。以语音识别，即将声学语音信号转换成 ASCII 文本为例。看上去我们可以毫无困难地做这件事，但是我们却不能解释我们是如何做的。由于年龄、性别或口音的差异，不同的人读相同的词发音却不同。在机器学习中，这个问题的解决方法是从不同的人那里收集大量发音样本，并学习将它们映射到词。

另一种需要学习的情况是要解决的问题随时间变化或依赖于特定的环境。我们希望有一个能够自动适应环境的通用系统，而不是为每个特定的环境编写一个不同的程序。以计算机网络上的包传递为例。最大化服务质量的、从源地到目的地的路径随网络流量的改变而改变。学习路由程序能够通过监视网络流量自动调整到最佳路径。另一个例子是智能用户界面，它能够自动适应用户的生物特征，即用户的口音、笔迹、工作习惯等。

机器学习在各个领域都有许多成功的应用：已经有了识别语音和笔迹的商用系统。零售商分析他们过去的销售数据，了解顾客行为，以便改善顾客关系管理。金融机构分析过去的交易，以便预测顾客的信用风险。机器人学习优化它们的行为，以便使用最少的资源来完成任务。在生物信息学方面，使用计算机不仅可以分析海量数据，而且还可以提取知识。这些只是我们（即你和我）将在本书讨论的应用的一部分。我们只能想象一下可使用机器学习实现的未来应用：可以在不同的路况、不同的天气条件下自己行驶的汽车，可以实时翻译外语的电话，可以在新环境（例如另一个星球的表面）航行的自动化机器人。机器学习的确是一个令人激动的研究领域！

本书讨论的许多方法都源于各种领域：统计学、模式识别、神经网络、人工智能、信号处理、控制和数据挖掘。过去，这些不同领域的研究遵循不同的途径，侧重点也不同。本书旨在把它们组合在一起，给出问题的统一处理并提供它们的解。

本书是一本入门教材，用于高年级本科生和研究生的机器学习课程，以及在业界工作、对这些方法的应用感兴趣的工程技术人员。预备知识是计算机程序设计、概率论、微积分和线性代数方面的课程。本书的目标是充分解释所有的学习算法，使得从本书给出的方程到计算机程序只是一小步。为了使这一任务更容易完成，对于某些情况，我们给出了算法的伪代码。

适当选取一些章节，本书可用作一学期的课程。再额外讨论一些研究论文的话，本书也可以用作两学期的课程，这时每章后的参考文献将很有用。

本书网页为 <http://www.cmpe.boun.edu.tr/~ethem/i2ml/>，我将在那里提供一些与本书有关的信息，如勘误表。我真诚地欢迎你将反馈意见发到我的邮箱：alpaydin@boun.edu.tr。

我非常喜欢写这本书，希望你能喜欢读它。

获得好想法的途径是与有才干的人一起工作，与他们一起工作也是一种乐趣。Boğaziçi 大学计算机工程系是一个极好的工作场所，在我写这本书时，我的同事们为我提供了我所需的所有支持。我也要感谢我过去和现在的学生，在他们身上，我实际检验了现在写进本书中的内容。

在写本书时，我得到了土耳其科学院青年科学家奖励计划的资助（EATÜBA-GEBİP/2001-1-1）。

我特别感谢 Michael Jordan。对于他多年来的支持和最近对本书的支持，我深表感谢。他针对本书大体组织和第 1 章所给出的建议在内容和形式上都大大改进了本书。Taner Bilgiç、Vladimir Cherkassky、Tom Dietterich、Fikret Gürgen、Olca Taner Yildiz 和 MIT 出版社的未留名审稿人也部分阅读了本书，并提供了非常宝贵的反馈。我希望他们在注意到我采纳了他们的建议但却没有特别致谢时，能够体会到我的感激之情。当然，书中的错误和不足应当由我个人负责。

我的父母信任我，我感谢他们永恒的爱和支持。无论我何时需要，Sema Oktuğ 总在身边，我将永远感激她的友谊。我还要感谢 Hakan Ünlü，在过去的几年中，我们无数次讨论了与生活、宇宙和万事万物相关的众多主题。

本书使用 Chris Manning 准备的 LATEX 宏排版，对此我很感谢他。我要感谢 MIT 出版社的编辑们，以及 Bob Prior、Valerie Geary、Kathleen Caruso、Sharon Deacon Warne、Erica Schultz 和 Emily Gutheinz，感谢他们在本书完成期间的不断支持和帮助。

自从第1版2004年问世以来，机器学习已经取得重要进展。首先，在应用领域迅速扩展。现在，互联网相关的技术，如搜索引擎、推荐系统、垃圾邮件过滤和入侵检测系统都在常规地使用机器学习。在生物信息学和计算生物学领域，由数据学习的方法被越来越广泛地使用。在自然语言处理应用（例如，机器翻译）中，我们看到越来越快地从编程的专家系统过渡到从实例文本的大型语料库中自动地学习。在机器人、医疗诊断、语音和图像识别、生物测定学、财经等领域，我们看到了本书讨论的机器学习方法越来越多的应用。这些应用有时是在模式识别的名义下，有时乔装为数据挖掘或披着某种其他外衣。

其次，理论上取得了重要进步。尤其是，核函数的思想和使用核函数的核机器使得我们可以更好地表示问题，并且与使用梯度下降训练的、具有S形隐藏单元的多层感知器相比，核机器的凸优化工作前进了一大步。贝叶斯方法通过选定适当的先验分布将专家的知识添加到数据隐含的知识上。图模型允许使用相互关联的节点的网络表示变量之间的依赖，并且有效地推断算法使得我们可以查询该网络。因此，有必要将核方法、贝叶斯估计和图模型这三个在第1版只用几节介绍的主题用更长的篇幅处理，扩充为新的三章。

该领域另一个极其重要的事情是认识到需要更好地设计机器学习实验。从使用单个检验集到交叉验证，再到配对 t 检验，我们已经走了很长的路。这就是为什么我要在第2版重写统计检验这一章，使之包含机器学习实验的设计与分析。要点是，检验不應該是在所有运行都完成之后再做的单独步骤（尽管介绍机器学习实验的新一章被安排在本书最后），应该提前设计实验的整个过程、定义相关因素、确定合适的实验过程。这时并且只有这时才做实验并分析实验结果。

长期以来，人们尤其是科学界的年长者相信，要想使机器像我们一样有智能，即要想使人工智能成为现实，无论就一般而言还是特殊地就计算机科学而言，我们当前的知识都是不够的。大多数人都认为，我们需要新的技术、新型材料、新型计算装置或新的程序设计技术，并且即便如此，我们也只能以有限的方式“模拟”人类智能的某些方面，而不可能完全实现人类的智能。

我相信我们很快将证明他们是错误的。我们最早在国际象棋中看到了这种证明，并且现在我们正在各个领域看到这种证明。有了足够的内存和计算能力，我们可以使用相对简单的算法来完成任务；这里的技巧是学习，或者从实例数据中学习，或者使用增强学习通过试错学习。看来，（例如，机器翻译）使用监督学习并且更多的是非监督学习算法很快将会成为可能。对于许多其他领域，例如使用增强学习的机器人无人导航也是如此。我相信在人工智能的许多领域这种态势都将继续，而关键是学习。只要我们为机器提供足够的数据（不必是监督的）和计算能力，如果机器可以自己学习，我们则不需要提出新的算法。

感谢世界各地使用第 1 版的所有教师和学生。感谢给我发送评论和勘误或以任何方式提供反馈的人。请继续给我发这样的邮件。我的 email 地址是 alpaydin@boun.edu.tr。

第 2 版还在网站上提供了更多的支持。本书的网站地址是 <http://www.cmpe.boun.edu.tr/~ethem/i2ml>。

感谢过去和现在跟我做论文的学生 Mehmet Gönen、Esma Kihç、Murat Semerci、M. Aydın Ulaş 和 Olcay Taner Yildiz，感谢在过去的几年选修 CmpE 544、CmpE 545、CmpE 591 和 CmpE 58E 课程的学生。检查你的知识的最佳方法是讲述它。

第 2 版再次与 MIT 出版社合作是一件令人愉快的事。感谢 Bob Prior、Ada Brunstein、Erin K. Shoudy、Kathleen Caruso 和 Marcy Ross 的帮助和支持。

符号表

Introduction to Machine Learning, Second Edition

x	标量值
\mathbf{x}	向量
\mathbf{X}	矩阵
\mathbf{x}^T	转置
\mathbf{X}^{-1}	逆矩阵
X	随机变量
$P(X)$	概率质量函数, X 是离散的
$p(X)$	概率密度函数, X 是连续的
$P(X Y)$	给定 Y , X 的条件概率
$E[X]$	随机变量 X 的期望值
$\text{Var}(X)$	X 的方差
$\text{Cov}(X, Y)$	X 和 Y 的协方差
$\text{Corr}(X, Y)$	X 和 Y 的相关性
μ	均值
σ^2	方差
Σ	协方差矩阵
m	均值的估计
s^2	方差的估计
S	协方差矩阵的估计
$\mathcal{N}(\mu, \sigma^2)$	一元正态分布, 均值为 μ , 方差为 σ^2
\mathcal{Z}	单位正态分布: $\mathcal{N}(0, 1)$
$\mathcal{N}_d(\mu, \Sigma)$	d -变量正态分布, 均值向量为 μ , 协方差矩阵为 Σ
x	输入
d	输入数 (输入的维度)
y	输出
r	要求的输出
K	输出数(类)
N	训练实例数
z	隐藏的值, 内蕴维, 潜在因子
k	隐藏维数, 潜在因子数
C_i	类 i
\mathcal{X}	训练样本

$\{x^t\}_{t=1}^N$	x 的集合, 上标 t 遍取 1 到 N
$\{x^t, r^t\}_t$	上标为 t 的输入和期望输出的有序对的集合
$g(x \theta)$	x 的函数, 其定义依赖于参数集 θ
$\arg \max_{\theta} g(x \theta)$	参数 θ , g 关于它取最大值
$\arg \min_{\theta} g(x \theta)$	参数 θ , g 关于它取最小值
$E(\theta \chi)$	样本 χ 上具有参数 θ 的误差函数
$l(\theta \chi)$	样本 χ 上具有参数 θ 的似然函数
$\mathcal{L}(\theta \chi)$	样本 χ 上具有参数 θ 的对数似然函数
$1(c)$	如果 c 为真, 则值为 1, 否则为 0
$\#\{c\}$	c 为真的元素数目
δ_{ij}	Kronecker δ : 如果 $i=j$, 取 1, 否则取 0

目 录

Introduction to Machine Learning, Second Edition

出版者的话	
中文版序	
译者序	
前言	
致谢	
关于第2版	
符号表	
第1章 绪论	1
1.1 什么是机器学习	1
1.2 机器学习的应用实例	3
1.2.1 学习关联性	3
1.2.2 分类	3
1.2.3 回归	6
1.2.4 非监督学习	7
1.2.5 增强学习	8
1.3 注释	8
1.4 相关资源	10
1.5 习题	11
1.6 参考文献	12
第2章 监督学习	13
2.1 由实例学习类	13
2.2 VC维	15
2.3 概率逼近正确学习	16
2.4 噪声	17
2.5 学习多类	18
2.6 回归	19
2.7 模型选择与泛化	21
2.8 监督机器学习算法的维	23
2.9 注释	24
2.10 习题	25
2.11 参考文献	25
第3章 贝叶斯决策定理	27
3.1 引言	27
3.2 分类	28
3.3 损失与风险	29
3.4 判别式函数	31
3.5 效用理论	31
3.6 关联规则	32
3.7 注释	33
3.8 习题	33
3.9 参考文献	34
第4章 参数方法	35
4.1 引言	35
4.2 最大似然估计	35
4.2.1 伯努利密度	36
4.2.2 多项密度	36
4.2.3 高斯(正态)密度	37
4.3 评价估计: 偏倚和方差	37
4.4 贝叶斯估计	38
4.5 参数分类	40
4.6 回归	43
4.7 调整模型的复杂度: 偏倚/方差 两难选择	45
4.8 模型选择过程	47
4.9 注释	50
4.10 习题	50
4.11 参考文献	51
第5章 多元方法	52
5.1 多元数据	52
5.2 参数估计	52
5.3 缺失值估计	53
5.4 多元正态分布	54
5.5 多元分类	56
5.6 调整复杂度	59
5.7 离散特征	61
5.8 多元回归	62
5.9 注释	63
5.10 习题	63
5.11 参考文献	64

第 6 章 维度归约	65	8.8 注释	110
6.1 引言	65	8.9 习题	111
6.2 子集选择	65	8.10 参考文献	112
6.3 主成分分析	67	第 9 章 决策树	113
6.4 因子分析	71	9.1 引言	113
6.5 多维定标	75	9.2 单变量树	114
6.6 线性判别分析	77	9.2.1 分类树	114
6.7 等距特征映射	80	9.2.2 回归树	118
6.8 局部线性嵌入	81	9.3 剪枝	119
6.9 注释	83	9.4 由决策树提取规则	120
6.10 习题	84	9.5 由数据学习规则	121
6.11 参考文献	85	9.6 多变量树	124
第 7 章 聚类	86	9.7 注释	125
7.1 引言	86	9.8 习题	126
7.2 混合密度	86	9.9 参考文献	127
7.3 k -均值聚类	87	第 10 章 线性判别式	129
7.4 期望最大化算法	90	10.1 引言	129
7.5 潜在变量混合模型	93	10.2 推广线性模型	130
7.6 聚类后的监督学习	94	10.3 线性判别式的几何意义	131
7.7 层次聚类	95	10.3.1 两类问题	131
7.8 选择簇个数	96	10.3.2 多类问题	132
7.9 注释	96	10.4 逐对分离	132
7.10 习题	97	10.5 参数判别式的进一步讨论	133
7.11 参考文献	97	10.6 梯度下降	135
第 8 章 非参数方法	99	10.7 逻辑斯谛判别式	135
8.1 引言	99	10.7.1 两类问题	135
8.2 非参数密度估计	99	10.7.2 多类问题	138
8.2.1 直方图估计	100	10.8 回归判别式	141
8.2.2 核估计	101	10.9 注释	142
8.2.3 k -最近邻估计	102	10.10 习题	143
8.3 到多元数据的推广	103	10.11 参考文献	143
8.4 非参数分类	104	第 11 章 多层感知器	144
8.5 精简的最近邻	105	11.1 引言	144
8.6 非参数回归: 光滑模型	106	11.1.1 理解人脑	144
8.6.1 移动均值光滑	106	11.1.2 神经网络作为并行 处理的典范	145
8.6.2 核光滑	108	11.2 感知器	146
8.6.3 移动线光滑	108	11.3 训练感知器	148
8.7 如何选择光滑参数	109		

11.4 学习布尔函数	150	12.11 习题	190
11.5 多层感知器	151	12.12 参考文献	190
11.6 作为普适近似的 MLP	153	第 13 章 核机器	192
11.7 后向传播算法	154	13.1 引言	192
11.7.1 非线性回归	154	13.2 最佳分离超平面	193
11.7.2 两类判别式	157	13.3 不可分情况: 软边缘超平面	195
11.7.3 多类判别式	158	13.4 ν -SVM	197
11.7.4 多个隐藏层	158	13.5 核技巧	198
11.8 训练过程	158	13.6 向量核	199
11.8.1 改善收敛性	158	13.7 定义核	200
11.8.2 过分训练	159	13.8 多核学习	201
11.8.3 构造网络	161	13.9 多类核机器	202
11.8.4 线索	162	13.10 用于回归的核机器	203
11.9 调整网络规模	163	13.11 一类核机器	206
11.10 学习的贝叶斯观点	164	13.12 核维度归约	208
11.11 维度归约	165	13.13 注释	209
11.12 学习时间	167	13.14 习题	209
11.12.1 时间延迟神经网络	167	13.15 参考文献	210
11.12.2 递归网络	168	第 14 章 贝叶斯估计	212
11.13 注释	169	14.1 引言	212
11.14 习题	170	14.2 分布参数的估计	213
11.15 参考文献	170	14.2.1 离散变量	213
第 12 章 局部模型	173	14.2.2 连续变量	215
12.1 引言	173	14.3 函数参数的贝叶斯估计	216
12.2 竞争学习	173	14.3.1 回归	216
12.2.1 在线 k -均值	173	14.3.2 基函数或核函数的使用	218
12.2.2 自适应共鸣理论	176	14.3.3 贝叶斯分类	219
12.2.3 自组织映射	177	14.4 高斯过程	221
12.3 径向基函数	178	14.5 注释	223
12.4 结合基于规则的知识	182	14.6 习题	224
12.5 规范化基函数	182	14.7 参考文献	224
12.6 竞争的基函数	184	第 15 章 隐马尔可夫模型	225
12.7 学习向量量化	186	15.1 引言	225
12.8 混合专家模型	186	15.2 离散马尔可夫过程	225
12.8.1 协同专家模型	188	15.3 隐马尔可夫模型	227
12.8.2 竞争专家模型	188	15.4 HMM 的三个基本问题	229
12.9 层次混合专家模型	189	15.5 估值问题	229
12.10 注释	189	15.6 寻找状态序列	231

15.7 学习模型参数	233	17.12 注释	270
15.8 连续观测	235	17.13 习题	271
15.9 带输入的 HMM	236	17.14 参考文献	272
15.10 HMM 中的模型选择	236	第 18 章 增强学习	275
15.11 注释	237	18.1 引言	275
15.12 习题	238	18.2 单状态情况: K 臂赌博机问题	276
15.13 参考文献	239	18.3 增强学习基础	277
第 16 章 图方法	240	18.4 基于模型的学习	278
16.1 引言	240	18.4.1 价值迭代	279
16.2 条件独立的典型情况	241	18.4.2 策略迭代	279
16.3 图模型实例	245	18.5 时间差分学习	280
16.3.1 朴素贝叶斯分类	245	18.5.1 探索策略	280
16.3.2 隐马尔可夫模型	246	18.5.2 确定性奖励和动作	280
16.3.3 线性回归	248	18.5.3 非确定性奖励和动作	282
16.4 d -分离	248	18.5.4 资格迹	283
16.5 信念传播	249	18.6 推广	285
16.5.1 链	249	18.7 部分可观测状态	286
16.5.2 树	250	18.7.1 场景	286
16.5.3 多树	251	18.7.2 例子: 老虎问题	287
16.5.4 结树	252	18.8 注释	290
16.6 无向图: 马尔可夫随机场	253	18.9 习题	291
16.7 学习图模型的结构	254	18.10 参考文献	292
16.8 影响图	255	第 19 章 机器学习实验的	
16.9 注释	255	设计与分析	294
16.10 习题	256	19.1 引言	294
16.11 参考文献	256	19.2 因素、响应和实验策略	296
第 17 章 组合多学习器	258	19.3 响应面设计	297
17.1 基本原理	258	19.4 随机化、重复和阻止	298
17.2 产生有差异的学习器	258	19.5 机器学习实验指南	298
17.3 模型组合方案	260	19.6 交叉验证和再抽样方法	300
17.4 投票法	261	19.6.1 K -折交叉验证	300
17.5 纠错输出码	263	19.6.2 5×2 交叉验证	301
17.6 装袋	265	19.6.3 自助法	302
17.7 提升	265	19.7 度量分类器的性能	302
17.8 重温混合专家模型	267	19.8 区间估计	304
17.9 层叠泛化	268	19.9 假设检验	307
17.10 调整系综	268	19.10 评估分类算法的性能	308
17.11 级联	269	19.10.1 二项检验	308

绪 论

1.1 什么是机器学习

为了在计算机上解决问题，我们需要算法。算法是指令的序列，它把输入变换成输出。例如，我们可以为排序设计一个算法，输入是数的集合，而输出是它们的有序列表。对于相同的任务可能存在不同的算法，而我们感兴趣的是如何找到需要的指令或内存最少，或者二者都最少的最有效算法。

然而，对于某些任务，我们没有算法；例如，我们没有将垃圾邮件与正常邮件分开的算法。我们知道输入是邮件文档，最简单的情况是一份字符文件。我还知道输出应该是指出消息是否为垃圾邮件的“是”或“否”，但是我们不知道如何把这种输入变换成输出。所谓的垃圾邮件随时间而变，因人而异。

我们缺乏的是知识，作为补偿我们有数据。我们可以很容易地编辑数以千计的实例消息，其中一些我们知道是垃圾邮件，而我们要做到的是希望从中“学习”垃圾邮件的结构。换言之，我们希望计算机(机器)自动地为这一任务提取算法。不需要学习如何将数排序，因为我们已经有这样的算法；但是，对于许多应用而言，我们确实没有算法，而是有实例数据。

随着计算机技术的发展，我们现在已经拥有存储和处理海量数据以及通过计算机网络从远程站点访问数据的能力。目前大多数的数据存取设备都是数字设备，记录的数据也很可靠。以一家连锁超市为例，它拥有遍布全国各地的数百家分店，并且在为数百万顾客提供数千种商品的零售服务。销售点的终端设备记录每笔交易的详细资料，包括日期、顾客识别码、购买商品和数量、消费总额等。这是典型的每日几个 G 字节的数据。连锁超市希望能够预测某种产品可能的顾客。对于这一任务，算法同样并非是显然的；它随时间而变，因地域而异。只有分析这些数据，并且将它转换为可以利用的信息时，这些存储的数据才能变得有用，例如做预测。

我们并不确切地知道哪些人倾向于购买这种口味的冰淇淋，或者这位作家的下一本书是什么，也不知道谁喜欢看这部新电影、访问这座城市，或点击这一链接。我们不能确切地知道哪些人比较倾向于购买哪种特定的商品，也不知道应该向喜欢读海明威作品的人推荐哪位作者。如果我们知道，我们就不需要任何数据分析；我们只管供货并记录下编码就可以了。但是，正因为我们不知道，所以才只能收集数据，并期望从数据中提取这些问题或相似问题的答案。

我们确信存在某种过程，可以解释我们所观测到的数据。尽管我们不清楚数据产生过程(例如顾客行为)的细节，但是，我们知道数据产生不是完全随机的。人们并不是去超市随机购买商品。当人们买啤酒时，也会买薯片；夏天买冰淇淋，而冬天则为 Glühwein[⊖]买香

⊖ Glühwein 是一种温热、有点甜味、加香料的葡萄酒。圣诞节期间，在欧洲很受欢迎。——译者注

料。数据中存在确定的模式。

我们也许不能够完全识别该过程，但是我们相信，我们能够构造一个好的并且有用的近似(good and useful approximation)。尽管这样的近似还不可能解释一切，但其仍然可以解释数据的某些部分。我们相信，尽管识别全部过程也许是不可能的，但是我们仍然能够发现某些模式或规律。这正是机器学习的定位。这些模式可以帮助我们理解该过程，或者我们可以使用这些模式进行预测：假定将来，至少是不远的将来，情况不会与收集样本数据时有很大的不同，则未来的预测也将有望是正确的。

机器学习方法在大型数据库中的应用称为数据挖掘(data mining)。类似的情况如大量的金属氧化物以及原料从矿山中开采出来，处理后产生少量非常珍贵的物质。同样地，在数据挖掘中，需要处理大量的数据以构建简单有用的模型，例如具有高精度的预测模型。数据挖掘的应用领域非常广泛：除零售业以外，在金融业，银行分析他们的历史数据，构建用于信用分析、诈骗检测、股票市场等方面的应用模型；在制造业，学习模型可以用于优化、控制以及故障检测等；在医学领域，学习程序可以用于医疗诊断等；在电信领域，通话模式的分析可用于网络优化和提高服务质量；在科学研究领域，比如物理学、天文学以及生物学的大量数据只有用计算机才可能得到足够快的分析。万维网(World Wide Web)是巨大的，并且在不断地增长，因此在万维网上检索相关信息不可能依靠人工完成。

然而，机器学习不仅仅是数据库方面的问题，它也是人工智能的组成部分。为了智能化，处于变化环境中的系统必须具备学习的能力。如果系统能够学习并且适应这些变化，那么系统的设计者就不必预见所有的情况，并为它们提供解决方案了。

机器学习还可以帮助我们解决视觉、语音识别以及机器人方面的许多问题。以人脸识别问题为例：我们做这件事毫不费力；即使姿势、光线、发型等不同，我们每天还是可以通过看真实的面孔或其照片来认出我们的家人和朋友。但是我们做这件事是下意识的，而且无法解释我们是如何做的。因为我们不能够解释我们所具备的这种技能，我们也就不可能编写相应的计算机程序。但是我们知道，脸部图像并非只是像素点的随机组合；人脸是有结构的、对称的。脸上有眼睛、鼻子和嘴巴，并且它们都位于脸的特定部位。每个人的脸都有各自的眼睛、鼻子和嘴巴的特定组合模式。通过分析一个人脸部图像的多个样本，学习程序可以捕捉到那个人特有的模式，然后在所给的图像中检测这种模式，从而进行辨认。这就是模式识别(pattern recognition)的一个例子。

机器学习使用实例数据或过去的经验训练计算机，以优化某种性能标准。我们有依赖于某些参数的模型，而学习就是执行计算机程序，利用训练数据或以往经验来优化该模型的参数的。模型可以是预测性的(predictive)，用于未来的预测，或者是描述性的(descriptive)，用于从数据中获取知识，也可以二者兼备。

机器学习在构建数学模型时利用了统计学理论，因为其核心任务就是从样本中推理。计算机科学的角色是双重的：第一，在训练时，我们需要求解优化问题以及存储和处理通常所面对的海量数据的高效算法。第二，一旦学习得到了一个模型，它的表示和用于推理的算法解也必须是高效的。在特定的应用中，学习或推理算法的效率，即它的空间复杂度和时间复杂度，可能与其预测精确度同样重要。

现在，让我们更详细地讨论一些应用领域的例子，以进一步深入了解机器学习的类型和用途。

1.2 机器学习的应用实例

1.2.1 学习关联性

在零售业,例如超市连锁店,机器学习的一个应用是购物篮分析(basket analysis)。它的任务是发现顾客所购商品之间的关联性:如果人们在购买商品 X 时也通常购买商品 Y ,而有一名顾客购买了商品 X 却没有购买商品 Y ,则他(或她)即是商品 Y 的潜在顾客。一旦我们发现这类顾客,我们就能针对他们实行打包销售策略。

为发现关联规则(association rule),我们对学习形如 $P(Y|X)$ 的条件概率感兴趣,其中 X 是我们知道的顾客已经购买的商品或商品集, Y 表示在条件 X 下可能购买的商品。

假定考察已有的数据,计算得到 $P(\text{chips}|\text{beer})=0.7$,那么我们就可以定义规则:

购买啤酒(beer)的顾客中有70%的人也买了薯片(chips)。

我们也许想要区分不同的顾客。针对这个问题,我们需要估计 $P(Y|X, D)$,其中 D 是顾客的一组属性,如性别、年龄、婚姻状况等,这里假定我们已经得到了这些属性信息。如果是考虑书店而不是超市的销售问题,则商品就可能是书或作者等。对于Web门户网站入口问题,项对应着到Web网页的链接,而我们可以估计用户可能点击的链接,并利用这些信息来预先下载这些网页,以取得更快的网页存取速度。

4

1.2.2 分类

信贷是金融机构(例如银行)借出的一笔钱,需要连本带息偿还,通常是分期偿还。对银行来说,重要的是能够提前预测贷款风险。这种风险是客户不履行义务和不全额还款的可能性。既要确保银行获利,又要确保不会因提供超出客户财力的贷款而给客户带来不便。

在信用评分(credit scoring)(Hand 1998)中,银行要计算在给定信贷额度和客户信息情况下的风险。客户信息包括我们已经获取的数据以及与计算客户财力相关的数据,即收入、存款、担保、职业、年龄、以往经济记录等。银行有以往贷款的记录,包括客户数据以及贷款是否偿还。通过这类特定的申请数据,我们可以推断出一般规则,表示客户属性及其风险性的关联性。也就是说,机器学习系统用一个模型来拟合过去的的数据,以便能够对新的申请计算风险,从而决定接受或拒绝该项申请。

这是一个分类(classification)问题的例子,这里有两个类:低风险客户和高风险客户。客户信息作为分类器的输入(input),分类器的任务是将输入指派到其中的一个类。

利用以往数据进行训练后,学习得到的规则可能具有如下形式

IF income $> \theta_1$ AND savings $> \theta_2$ THEN low-risk ELSE high-risk

其中 θ_1 和 θ_2 是合适的值(参见图1-1)。这是判别式(discriminant)的一个例子,它是将不同类的样本分开的函数。

有了这样的规则,其主要用途就是预测(prediction):一旦我们拥有拟合以往数据的规则,如果未来与过去类似,那么我们就能够对新的实例做出正确的预测。如果给定一个新的具有特定收入(income)和存款(savings)的申请,我们就可以很容易地判断出它是低风险

(low-risk) 还是高风险 (high-risk) 了。

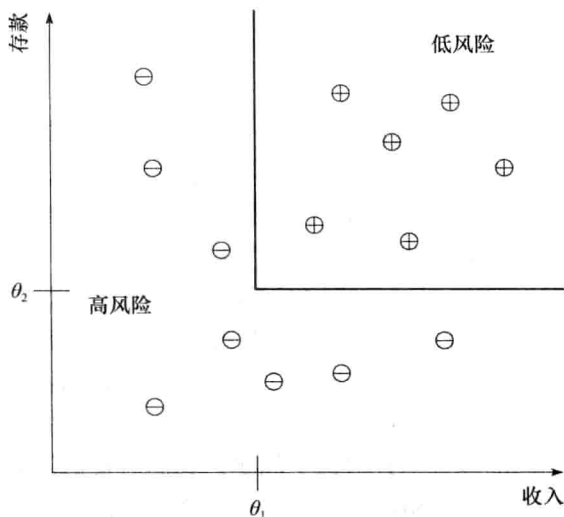


图 1-1 训练数据集示例，其中每个圆圈对应一个数据实例，输入值在对应的坐标上，符号则指示着类别。为简单起见，输入只包括客户的收入 (income) 和存款 (savings) 两种属性，两个类分别为低风险 (“+”) 和高风险 (“-”)。图中还显示了分隔两类样本的判别式样例

在某些情况下，我们可能不希望做 0/1 (低风险/高风险) 类型的判断，而是希望计算一个概率值 $P(Y|X)$ ，其中 X 是顾客属性， Y 是 0 或 1，分别表示低风险和高风险。从这个角度来看，我们可以将分类看作学习从 X 到 Y 的关联性。于是，给定 $X=x$ ，如果有 $P(Y=1|X=x)=0.8$ ，则我们就说该客户为高风险的可能性有 80%，或者等价地说该客户为低风险的可能性有 20%。然后，我们可以根据可能的收益和损失来决定接受或拒绝这笔贷款业务。

机器学习在模式识别 (pattern recognition) 方面有很多的应用。其中之一是光学字符识别 (optical character recognition, OCR)，即从字符图像识别字符编码。这是一个多类问题的例子，类与我们想要识别的字符一样多。特别有趣的是手写体字符的识别问题。人们有不同的书写风格；字体有大有小，倾斜角度不同，还有用钢笔或用铅笔之别，所以同一个字符可能会有许多种可能的图像。尽管书写是人类的发明创造，但是我们还没有像人类读者一样准确的系统。我们没有字符“A”的形式化描述，涵盖所有“A”而不涵盖任何非“A”。没有这种形式化描述，我们就要从书写者那里取样，从这些实例中学习关于“A”的定义。然而，尽管我们不知道是什么因素使得一个图像被识别为“A”，但是我们确信所有这些不同的“A”的图像都具有某些共同的特征，这正是我们希望从实例中提取的。我们知道，图像不只是随机点的集合，它是笔画的集合，并且是有规律的，通过学习程序我们能够捕获这些规律。

阅读文本时，我们能够利用的一个因素是人类语言的冗余性。词是字符的序列 (sequence)，并且相继的符号不是独立的，而是被语言的词所约束。这有好处，即便有一个符号不能识别，我们仍可以读出词 t?e[⊖]。根据语言的语法和语义，这种上下文的依赖性还

⊖ 这里，“?”表示不能识别的符号。——译者注

可能出现在词和句子之间等较高的层次上。目前有用于学习序列和对这种依赖性建模的机器学习算法。

对于人脸识别(face recognition), 输入是人脸图像, 类是需要识别的人, 并且学习程序应当学习人脸图像与身份之间的关联性。这个问题比 OCR 更困难, 原因是人脸会有更多的类, 输入图像也更大一些, 并且人脸是三维的, 不同的姿势和光线等都会导致图像的显著变化。另外, 对于特定人脸的输入也会出现问题, 比如说眼镜可能会把眼睛和眉毛遮住, 胡子可能会把下巴盖住等。

对于医学诊断(medical diagnosis), 输入是关于患者的信息, 而类是疾病。输入包括患者的年龄、性别、既往病史、目前症状等。当然, 患者可能还没有做过某些检查, 因此这些输入将会缺失。检查需要时间, 还可能要花很多钱, 而且也许还会给患者带来不便。因此, 除非我们确信检查将提供有价值的信息, 否则我们将不对患者进行检查。在医学诊断的情况下, 错误的诊断结果可能会导致我们采取错误的治疗或根本不进行治疗。在不能确信诊断结果的情况下, 分类器最好还是放弃判定, 而等待医学专家来做决断。

在语音识别(speech recognition), 输入是语音, 类是可以读出的词汇。这里要学习的是从语音信号到某种语言的词汇的关联性。由于年龄、性别或口音方面的差异, 不同的人对于相同词汇的读音不同, 这使得语音识别问题相当困难。语音识别的另一个特点是其输入信号是时态的(temporal), 词汇作为音素的序列实时读出, 而且有些词汇的读音会较长一些。

语音信息的作用有限, 并且与光学字符识别一样, 在语音识别中, “语言模型”的集成是至关重要的, 而且提供语言模型的最好方法仍然是从实例数据的大型语料库中学习。机器学习在自然语言处理(natural language processing)方面的应用与日俱增。垃圾邮件过滤就是一种应用, 其中垃圾邮件的制造者为一方, 过滤者为另一方, 一直都在寻找越来越精巧的方法以便超越对方, 也许最吸引人的是机器翻译(machine translation)。经历了数十年手工编写翻译规则的研究之后, 最近人们认识到最有希望的方法是提供大量翻译文本实例对, 并且让程序自动地提取一个字符串映射到另一个字符串的规则。

生物测定学(biometrics)使用人的生理和行为特征来识别或认证人的身份, 需要集成来自不同形态的输入。生理特征的例子是面部图像、指纹、虹膜和手掌; 行为特征的例子是签字的力度、嗓音、步态和击键。与通常的鉴别过程(照片、印刷签名或口令)相反, 会有许多不同的(不相关的)输入, 伪造(欺骗)更困难, 并且系统更准确, 不会对用户太不方便。机器学习既用于针对这些不同形态而构建不同的识别器, 也考虑这些不同数据源的可靠性, 用于组合它们的决策, 以便得到接受或拒绝的总体决策。

从数据中学习规则也为知识抽取(knowledge extraction)提供了可能性。规则是一种解释数据的简单模型, 而观察该模型我们能得到潜在数据处理的解释。例如, 一旦我们学会了区分低风险客户和高风险客户的判别式, 我们就拥有了关于低风险客户特性的知识。然后, 我们就能够利用这些知识, 通过比如广告等方式, 更有效地争取那些潜在的低风险客户。

机器学习还可以进行压缩(compression)。用规则拟合数据, 我们能得到比数据更简单的解释, 需要的存储空间更少, 处理所需要的计算更少。例如, 一旦你掌握了加法规则, 你就不必记忆每对可能数字的和是多少。

机器学习的另一种用途是离群点检测(outlier detection), 即发现那些不遵守规则的例外

实例。在这种情况下，学习规则之后，我们感兴趣的不是规则，而是规则未能覆盖的例外，它们可能暗示出我们需要注意的异常，如诈骗。

1.2.3 回归

假设我们想要一个能够预测二手车价格的系统。该系统的输入是我们认为会影响到车价的属性信息：品牌、车龄、发动机性能、里程以及其他信息。输出是车的价格。这种输出为数值的问题是回归 (regression) 问题。

设 x 表示车的属性， y 表示车的价格。调查一下以往的交易情况，我们能够收集训练数据，而机器学习程序用一个函数拟合这些数据来学习 x 的函数 y 。图 1-2 给出了一个例子，其中对于 w 和 w_0 的合适值，拟合函数具有以下形式：

$$y = wx + w_0$$

回归和分类均为监督学习 (supervised learning) 问题，其中输入 x 和输出 y 给定，任务是学习从输入到输出的映射。机器学习的方法是，先假定某个依赖于一组参数的模型：

$$y = g(x | \theta)$$

其中， $g(\cdot)$ 是模型， θ 是模型的参数。对于回归， y 是数值；对于分类， y 是类编码 (如 0/1)。 $g(\cdot)$ 为回归函数，或者 (对于分类) 是将不同类的实例分开的判别式函数。机器学习程序优化参数 θ ，使得逼近误差最小，也就是说，我们的估计要尽可能地接近训练集中给定的正确值。例如，图 1-2 所示的模型是线性的， w 和 w_0 是为最佳拟合训练数据优化的参数。在线性模型限制过强的情况下，我们可以利用比如二次函数：

$$y = w_2 x^2 + w_1 x + w_0$$

或更高阶的多项式，或其他非线性函数，为最佳拟合优化它们的参数。

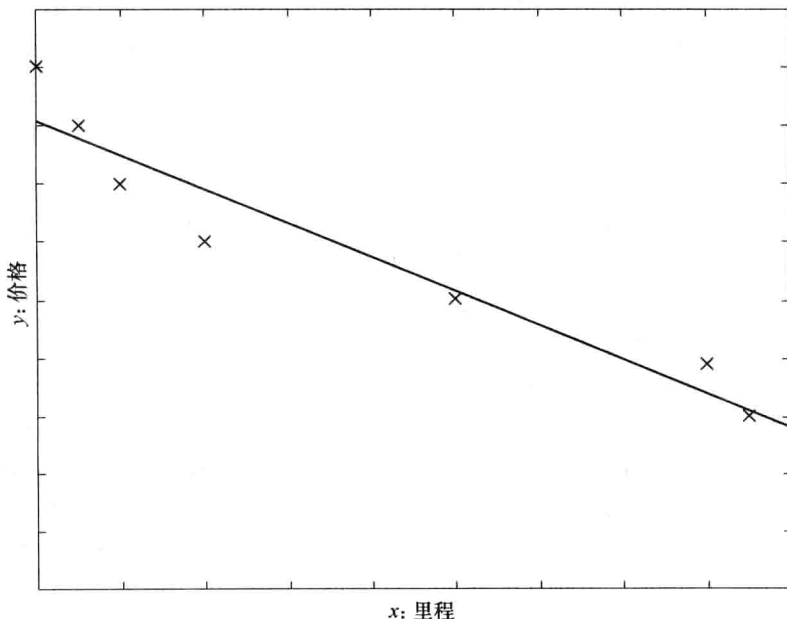


图 1-2 二手车的训练数据及其拟合函数。为简单起见，这里采用线性模型，输入属性也只有里程

回归的另一个例子是对移动机器人的导航。例如，自动汽车导航。其中输出是每次转动车轮的角度，使得汽车前进而不会撞到障碍物或偏离车道。这种情况下，输入由汽车上的传感器（如视频相机、GPS 等）提供。训练数据可以通过监视和记录驾驶员的动作收集。

我们来想象回归的其他应用，这里我们试图优化一个函数^①。假设我们想要造一个焙炒咖啡的机器，该机器有多个影响咖啡品质的输入：各种温度、时间、咖啡豆种类等配置。我们针对不同的输入配置进行大量试验，并测量咖啡的品质，例如，根据消费者的满意度测量咖啡的品质。为寻求最优配置，我们拟合一个联系这些输入和咖啡品质的回归模型，并在当前模型的最优样本附近选择一些新的点，以便寻找更好的配置。我们抽取这些点，检测咖啡的品质，将它们加入训练数据，并拟合新的模型。这通常被称为响应面设计（response surface design）。

10

1.2.4 非监督学习

在监督学习中，我们的目标是学习从输入到输出的映射关系，其中输出的正确值已经由指导者提供。然而，非监督学习中却没有这样的指导者，只有输入数据。我们的目标是发现输入数据中的规律。输入空间存在着某种结构，使得特定的模式比其他模式更常出现，而我们希望知道哪些经常发生，哪些不经常发生。在统计学中，这称为密度估计（density estimation）。

密度估计的一种方法是聚类（clustering），其目标是发现输入数据的簇或分组。对于拥有老客户数据的公司，客户数据包括客户的个人统计信息，及其以前与公司的交易，公司也许想知道其客户的分布，搞清楚什么类型的客户会频繁出现。这种情况下，聚类模型会将属性相似的客户分派到相同的分组，为公司提供其客户的自然分组；这称作客户市场划分（customer segmentation）。一旦找出了这样的分组，公司就可能做出一些决策，比如对不同分组的客户提供特别的服务和产品等；这称作客户关系管理（customer relationship management）。这样的分组也可以用于识别“离群点”，即那些不同于其他客户的客户，这可能意味着一块新的市场，公司可以进一步开发。

聚类的一个有趣的应用是图像压缩（image compression）。在这种情况下，输入实例是由 RGB 值表示的图像像素。聚类程序将颜色近似的像素点分到相同的分组，而这样的分组对应图像中频繁出现的颜色。如果图像中只有少数颜色，并且我们用一种颜色（例如颜色的平均值）对属于同一分组的像素进行编码，则图像被量化。假设像素是 24 位，表示 1 600 万种颜色，但是如果只有 64 种主色调，那么对于每个像素，我们只需要 6 位而不是 24 位。例如，如果景象在图像的不同部分有多种不同的蓝色色调，并且我们采用它们的平均值来表示所有这些蓝色，那么我们就丢失了图像的细节，但是赢得了存储和传送图像的空间。理想状态下，人们希望通过分析重复的图像模式（如纹理、对象等）来识别更高层次的规律性。这为更高层次、更简单、更有用地描述景象提供了可能，并且实现了比像素级更好的压缩。如果我们扫描了文档页，我们得到的不是一些随机的有/无像素，而是一些字符的位图。这样的数据是有结构的，并且我们利用这些冗余信息找出数据的较短描述：“A”的 16 × 16 的位图占 32 字节，其 ASCII 码只占 1 个字节。

11

① 感谢 Michael Jordan 提供这个例子。

在文档聚类(document clustering)中,目标是把相似的文档分组。例如,新闻报道可以进一步划分为涉及政治、体育、时尚、艺术等子组。通常,文档用词袋(bag of words)表示;即预先定义 N 个词的词典,并且每份文档都是一个 N 维二值向量,如果第 i 个词出现在该文档中,则其第 i 个分量取1。删除后缀“-s”和“-ing”等,以避免重复,并且不用诸如“of”、“and”等不包含信息的词。然后,文档根据它们包含的相同词的个数分组。当然,如何选取词典是至关重要的。

机器学习方法还应用于生物信息学(bioinformatics)。在我们的基因组中,DNA是“生命的蓝图”,也是碱基即A、G、C和T的序列。RNA由DNA转录而来,蛋白质由RNA转换而来。蛋白质就是生命体和生命体的产物。正如DNA是碱基序列,蛋白质则是氨基酸(由碱基定义)序列。计算机科学在分子生物学的应用领域之一就是比对(alignment),即将一个序列与另一个序列匹配。这是一个困难的串匹配问题,因为序列可能相当长,有很多模板串要进行匹配,并且还可能会被删节、插入和置换。聚类用于学习结构域(motif),这是蛋白质结构中反复出现的氨基酸序列。结构域之所以令人感兴趣,是因为它们可能对应它们所表征的序列内部的结构或功能要素。比方说,如果氨基酸是字母,蛋白质是句子,那么结构域就像单词,即具有特别意义、频繁地出现在不同句子中的一串字母。

[12]

1.2.5 增强学习

在某些应用中,系统的输出是动作(action)的序列。在这种情况下,单个的动作并不重要,重要的是策略(policy),即达到目标的正确动作的序列。不存在中间状态中最好动作这种概念。如果一个动作是好的策略的组成部分,那么该动作就是好的。这种情况下,机器学习程序就应当能够评估策略的好坏程度,并从以往好的动作序列中学习,以便能够产生策略。这种学习方法称为增强学习(reinforcement learning)算法。

游戏(game playing)是一个很好的例子。在游戏中,单个移动本身并不重要,正确的移动序列才是重要的。如果一个移动是一个好的游戏策略的一部分,则它就是好的。游戏是人工智能和机器学习的重要研究领域,这是因为游戏容易描述,但又很难玩好。像国际象棋这样的游戏,其规则只有少量的几条,但是它非常复杂,因为在每种状态下都有大量可行的移动,并且每局又都包含有大量的移动。一旦有了能够学习如何玩好游戏的好算法,我们也可以将这些算法用在具有更显著经济效益的领域。

用于在某种环境下搜寻目标位置的机器人导航是增强学习的另一个应用领域。在任何时候,机器人都能够朝着多个方向之一移动。经过多次的试运行,机器人应当学到正确的动作序列,尽可能快地从某一初始状态到达目标状态,并且不会撞到任何障碍物。致使增强学习难度增加的一个因素是系统具有不可靠和不完整的感知信息。例如,装备视频照相机的机器人就得不到完整的信息,因此该机器人总是处于部分可观测(partially observable)状态,并且应当将这种不确定性考虑在内。一个任务还可能需要多智能主体(multiple agents)的并行操作,这些智能主体将相互作用并协同操作,以便完成一个共同的目标。机器人足球是这种情况的例子之一。

[13]

1.3 注释

进化是形成我们的身体形状和我们的内在本能的主要力量。我们还需要终生学习,以改

变我们的行为。这有助于我们适应进化论还不能预测的环境变化。在合适的环境下,具有短暂寿命的生物体可能具备它们所有天生的行为能力,但是上苍并未赋予我们应对在有限生命中可能遇见的所有状况的能力。但是,进化赋予我们大脑和学习机制,使得我们可以根据经验实现自我更新,从而适应各种各样的环境。当我们在特定情境下学习到最好的策略时,知识就存储在我们的大脑里。当情境再现,当我们再认知(“认知”意味认出)情境时,我们就能够回忆起合适的策略并采取相应的动作。不过,学习有其局限性;就我们大脑有限的容量来说,也许有些东西我们永远都不可能学会,正像我们永远不可能“学会”长出第三只手臂,或是在脑袋后面长只眼睛,即使它们是有用的我们也学不会。关于心理学视角下的学习和认知可参看 Leahey 和 Harris 1997。注意,与心理学、认知科学以及神经系统科学都不同,机器学习的目标并不是理解人类和动物学习的过程,而是像任何的工程领域一样,机器学习旨在构建一个有用的系统。

几乎所有的科学领域都在用模型拟合数据。科学家们设计实验、进行观测并收集数据。然后,通过找寻能解释所观测数据的简单模型,尝试抽取知识。该过程称为归纳(induction),它是从一组特别的示例中提取通用规则的过程。

现在,这样的数据分析已经不能再依赖人工完成了,原因有二:一是数据量巨大,二是能够做这种分析的人非常少而且人工分析又很昂贵。因而,对于能够分析数据并自动从中提取信息的计算机模型,也就是说对于学习,人们的兴趣正在不断地增长。

下面章节中我们将要讨论的方法源于不同的科学领域。有时,相同的算法会在多个领域

14

中沿着各自不同的历史轨迹被独立地发现。在统计学中,从特殊的观测到一般的描述称为推断(inference),而学习称为估计(estimation)。分类在统计学中称为判别式分析(discriminant analysis)(McLachlan 1992; Hastie、Tibshirani 和 Friedman 2001)。在计算机价格低廉并且数量充足之前,统计学家只能处理小样本。作为数学家,统计学家主要使用能够精确分析的简单参数模型。在工程学中,分类称为模式识别(pattern recognition),方法是非参数的,并且更大程度是凭借经验的(Duda、Hart 和 Stork 2001; Webb 1999)。机器学习与人工智能(artificial intelligence)有关(Russell 和 Norvig 1995),原因是智能系统应当能够适应其环境的变化。视觉、语音以及机器人等应用领域都是从样本数据中学习。在电子工程领域,信号处理(signal processing)的研究使得自适应计算机视觉和语音程序出现。其中,隐马尔可夫模型(Hidden Markov Models, HMM)的发展对于语音识别尤其重要。

20 世纪 80 年代后期,随着 VLSI 技术的发展和制造含有数千个处理器的并行硬件的可能性的出现,基于多处理单元的分布式计算理论的可行性使得人工神经网络(artificial neural networks)研究领域得到了重生(Bishop, 1995)。随着时间的推移,人们认识到在神经网络研究领域中,大多数的神经网络学习算法都具有其统计学的基础(例如多层感知器就是另一类的非参估计),因此模拟人脑计算的说法开始逐渐淡出。

近年来,基于核的算法(如支持向量机)日趋流行。借助于使用核函数,支持向量机适用于各种应用,尤其适合生物信息学和自然语言处理方面的应用。如今,人们已经广泛认识到,对于学习而言,好的数据表示至关重要,而核函数则是一种引进这种专家知识的好方法。

最近,随着存储和连接费用的降低,在因特网上使用非常大的数据集已经成为可能,再加上廉价的计算,已经使得在大量数据上运行学习算法成为可能。在过去的几十年中,人们一般相信,对于人工智能而言,我们需要新的范型、新的思维、新的计算模型或一些全新的算法。考虑到机器学习最近在各领域的成功,也许可以说,我们需要的不是新算法,而是大量数据实例和在这些数据上运行算法的充足计算能力。例如,支持向量机源于势函数(potential function),线性分类和基于最近邻的方法都是20世纪50或60年代提出的;那时,我们只是没有适合这些算法的快速计算机或大型存储器,才不能完全展示它们的潜力。可以推测,诸如机器翻译,甚至规划这样的任务都可以用这种相对简单的算法来解决,但需要在大量实例数据上训练,或通过长时间试错运行。智能看来不是源于某些稀奇古怪的公式,而是源于简单、直截了当的算法的耐心以及近乎蛮力的使用。

数据挖掘(data mining)的命名来源于机器学习算法在商界海量数据上的应用(Weiss 和 Indurkha 1998)。在计算机科学领域,数据挖掘也称为数据库知识发现(knowledge discovery in databases, KDD)。

在统计学、模式识别、神经网络信号处理、控制、人工智能以及数据挖掘等不同领域中,研究工作遵循着各自的途径,并有其各自的侧重点。本书的目标是结合所有这些研究重点,以给出统一的处理问题方法,并提出求解方案。

1.4 相关资源

机器学习的最新研究成果会发表在不同领域的会议和期刊上。机器学习专门的期刊有 *Machine Learning* (机器学习) 和 *Journal of Machine Learning Research* (机器学习研究)。以神经网络为主的期刊有 *Neural Computation* (神经计算)、*Neural Networks* (神经网络) 以及 *IEEE Transactions on Neural Networks* (IEEE 神经网络汇刊)。统计学方面的期刊如 *Annals of Statistics* (统计学年鉴) 和 *Journal of the American Statistical Association* (美国统计学会杂志) 也会发表一些机器学习方面的文章。另外, *IEEE Transactions on Pattern Analysis and Machine Intelligence* (IEEE 模式分析与机器智能汇刊) 也是机器学习研究性文章的资源之一。

关于人工智能、模式识别、模糊逻辑以及信号处理方面的期刊也包含机器学习方面的文章。以数据挖掘为主的期刊有 *Data Mining and Knowledge Discovery* (数据挖掘与知识发现)、*IEEE Transactions on Knowledge and Data Engineering* (IEEE 知识与数据工程汇刊) 以及 *ACM Special Interest Group on Knowledge Discovery and Data Mining Explorations Journal* (ACM 知识发现和数据挖掘特别兴趣组期刊)。

关于机器学习方面的主要会议有“*Neural Information Processing Systems*”(NIPS)、“*Uncertainty in Artificial Intelligence*”(UAI)、“*International Conference on Machine Learning*”(ICML)、“*European Conference on Machine Learning*”(ECML) 以及“*Computational Learning Theory*”(COLT)、“*International Joint Conference on Artificial Intelligence*”(IJCAI) 等。另外,关于神经网络、模式识别、模糊逻辑以及遗传算法等方面的会议,以及关于诸如计算机视觉、语音技术、机器人和数据挖掘等应用方面的会议,也会有针对机器学习的专题。

因特网上有很多数据集,致力于机器学习的研究者经常把它们作为自己研究的基准。下面是一些常用的网址:

- UCI 机器学习数据库是最流行的数据库: <http://www.ics.uci.edu/~mlearn/MLRepository.html>
- UCI KDD Archive: <http://kdd.ics.uci.edu/summary.data.application.html>
- Statlib: <http://lib.stat.cmu.edu>
- Delve: <http://www.cs.utoronto.ca/~delve/>

此外,还有一些针对特定应用程序的库,例如,针对计算生物学、人脸识别、语音识别等。

新的、更大的数据集不断地添加到这些数据库中,特别是添加到 UCI 机器学习库中。但是,一些研究者仍然相信这些库的范围有限,不能反映实际数据的全部特征,因此在这些库中数据集上的准确性并不说明问题。甚至可以说,当反复使用固定库中的数据集,量身打造新算法时,我们正在产生针对这些数据集的一组新的“UCI 算法”。

正如我们将在后面的章节中看到的,不同的算法在不同的任务上运行得更好,因此最好是针对一种应用,为该应用抽取一个或一些大型数据集,并针对特定的任务,在这些数据集上进行算法比较。

机器学习研究者近期的文章大多数都可以从互联网上找到, <http://citeseer.ist.psu.edu> 上的 NEC 研究索引是一个不错的研究入口,大部分作者还在网站上提供了他们的算法编码。还有一些实现各种机器学习算法的免费软件包,其中 Weka 特别值得关注: <http://www.cs.waikato.ac.nz/ml/weka>。

17

1.5 习题

1. 设想你有两种选择:可以传真一份文档,即传送图像;或者先使用光学字符阅读器(OCR),然后再传送相应的文本文件。用对比方式论述这两种方法的优缺点,并讨论什么时候一种方法比另一种方法更好。
2. 假定我们正在构建一个 OCR,并且对于每一个字符,我们都存储该字符的位图作为与逐个像素读取的字符匹配的模板。请解释什么时候这样的系统会失败?为什么条码读取器目前仍在使用?
3. 假定我们的既定目标是构建识别垃圾邮件的系统。请问是垃圾邮件中的什么特征使得我们能够确认它为垃圾邮件?计算机如何通过语法分析来发现垃圾邮件?如果发现了垃圾邮件,你希望计算机如何处理它:自动删除?转到另一个文件夹?还是仅仅在屏幕上标亮显示?
4. 如果给定任务是制造自动出租车,请定义约束。输入是什么?输出是什么?如何与乘客沟通?需要与其他的自动出租车沟通,即需要某种语言吗?
5. 在购物篮分析中,我们希望找出产品 X 和 Y 二者之间的依赖关系。对于给定的顾客交易数据库,如何能够发现这些数据之间的依赖关系?进而,如何将依赖关系发现算法推广到多于两个的产品之间?
6. 怎样能够预测用户下一次将键入的命令?或者,怎样能够预测 Web 上下一个将要被下载的网页?这样的预测什么时候是有用的?什么时候会变得令人讨厌?
7. 在你的日报中,为政治、体育和艺术类各找出 5 个新闻报道样例。阅读这些报道,找出

每类报道频繁使用的词，这些词可能帮助我们区别不同的类别。例如，政治方面的新闻报道多半会包含“政府”、“经济衰退”、“国会”等词，而在艺术类的新闻报道中可能包括“专辑”、“油画”或“剧院”。还有一些词(如“目标”)是模棱两可的。

- 18 8. 如果面部图像是 100×100 的图像，按行写出，则它是一个 10 000 维向量。如果我们把图像向右移动一个像素，则将得到 10 000 维空间中一个很不同的向量。如何构造一个对于这种扰动具有鲁棒性面部识别器？
9. 取一个词，例如“machine”。写 10 次，请一位朋友也写 10 次。分析这 20 个图像，试找出区分你与朋友手书的特征、笔画类型、曲度、圆和如何画点等。
10. 在估计二手车的价格时，估计它相对于原价的折旧率，而不是估计它的绝对价格则更有意义。为什么？

1.6 参考文献

- Bishop, C. M. 1995. *Neural Networks for Pattern Recognition*. Oxford: Oxford University Press.
- Duda, R. O., P. E. Hart, and D. G. Stork. 2001. *Pattern Classification*, 2nd ed. New York: Wiley.
- Han, J., and M. Kamber. 2006. *Data Mining: Concepts and Techniques*, 2nd ed. San Francisco: Morgan Kaufmann.
- Hand, D. J. 1998. “Consumer Credit and Statistics.” In *Statistics in Finance*, ed. D. J. Hand and S. D. Jacka, 69–81. London: Arnold.
- Hastie, T., R. Tibshirani, and J. Friedman. 2001. *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. New York: Springer.
- Leahey, T. H., and R. J. Harris. 1997. *Learning and Cognition*, 4th ed. New York: Prentice Hall.
- McLachlan, G. J. 1992. *Discriminant Analysis and Statistical Pattern Recognition*. New York: Wiley.
- Russell, S., and P. Norvig. 2002. *Artificial Intelligence: A Modern Approach*, 2nd ed. New York: Prentice Hall.
- Webb, A. 1999. *Statistical Pattern Recognition*. London: Arnold.
- Witten, I. H., and E. Frank. 2005. *Data Mining: Practical Machine Learning Tools and Techniques*, 2nd ed. San Francisco: Morgan Kaufmann.
- 19

监督学习

我们从最简单的情况开始讨论监督学习，首先从正例和负例集合中学习类别，继而推广并讨论多类的情况，然后再讨论输出为连续值的回归。

2.1 由实例学习类

假设我们要学习“家用汽车”类 C 。现在有一组汽车实例和一组被测人，我们向被测人示以这些汽车。被测人看到所示汽车并标记汽车，将他们认为家用汽车的标为正例 (positive example)，其他的标为负例 (negative example)。类学习就是找寻一个涵盖所有的正例而不包括任何负例的描述。通过这些，我们可以做预测：给定一辆我们以前从未见过的汽车，检查学习得到的描述，我们就可以判断这辆汽车是否为家用汽车。我们还可以做知识提取：这种研究可能由汽车公司赞助，目的可以是了解人们对家用汽车的期望。

经过与该领域专家的一些沟通，假定我们得到了一个结论：在我们所掌握的汽车的所有特征中，区别家用汽车与其他汽车的特征是价格和发动机功率。这两个属性就是类识别器的输入 (input)。注意，当我们决定采用这种特殊输入表示 (input representation) 时，我们忽略其他属性，将它们看作是不相关的。尽管有人可能认为诸如座位数量、车身颜色等属性对于辨别车型也很重要，但是这里为了简单起见，我们只考虑价格和发动机功率。

我们假设价格为第一个输入属性 x_1 (比如以美元计算)，发动机功率为第二个输入属性 x_2 (比如以立方厘米计发动机排量)。这样，每辆汽车就可以用两个数值来表示

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \quad (2.1)$$

而它的标号表示汽车的类型

$$r = \begin{cases} 1 & \text{如果 } \mathbf{x} \text{ 是正例} \\ 0 & \text{如果 } \mathbf{x} \text{ 是负例} \end{cases} \quad (2.2)$$

每辆汽车用一个这种有序对 (\mathbf{x}, r) 来表示，而训练集中包括 N 个这样的实例

$$\mathcal{X} = \{\mathbf{x}^t, r^t\}_{t=1}^N \quad (2.3)$$

其中， t 用于标记训练集中的各个汽车实例，它不表示时间或任何类似的序。

现在，我们的训练数据可以绘制在二维空间 (x_1, x_2) 上，其中每个实例 t 是一个数据点，坐标为 (x_1^t, x_2^t) ，其类型 (即正或负) 由 r^t 给定 (参见图 2-1)。

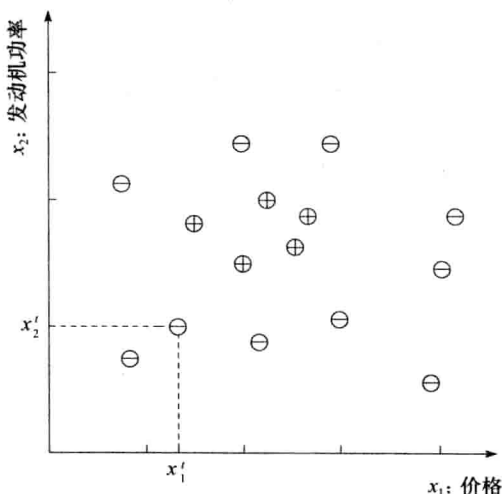


图 2-1 “家用汽车”类的训练集。其中每个点代表一个汽车实例，点的坐标值分别表示汽车的价格和发动机功率。“+”表示正例 (家用汽车)，“-”表示负例 (非家用汽车)，即其他类型的汽车

通过进一步与专家讨论和分析数据, 我们有理由相信, 对于家用汽车, 其价格和发动机功率应当是在某个确定的范围内

$$(p_1 \leq \text{价格} \leq p_2) \text{AND} (e_1 \leq \text{发动机功率} \leq e_2) \quad (2.4)$$

其中 p_1 , p_2 , e_1 和 e_2 为适当的值。这样, 式 (2.4) 假定类 C 是价格-发动机功率空间中的矩形 (参见图 2-2)。

式 (2.4) 确定了假设类 (hypothesis class) \mathcal{H} (即矩形的集合), 我们相信 C 是从中抽取的。学习算法应当找到一个特定的假设 (hypothesis) $h \in \mathcal{H}$, 尽可能地逼近 C 。

尽管专家定义了假设类, 但是他却不能说出参数值是什么。换句话说, 尽管我们选定了 \mathcal{H} , 但是我们却不知道哪个特定的 $h \in \mathcal{H}$ 等于或最接近 C 。然而, 一旦我们把注意力局限于这个假设类, 学习类就归结为较简单的问题——找出定义 h 的 4 个参数。

我们的目标是找出 $h \in \mathcal{H}$, 它与 C 尽可能类似。假设 h 对实例 x 进行预测, 使得

$$h(x) = \begin{cases} 1 & \text{如果 } h \text{ 将 } x \text{ 分类为正例} \\ 0 & \text{如果 } h \text{ 将 } x \text{ 分类为负例} \end{cases} \quad (2.5)$$

实际上我们并不知道 $C(x)$, 因此也无法评估 $h(x)$ 与 $C(x)$ 的匹配程度。我们所拥有的是训练集 X , 它是所有可能的 x 的一个小子集。经验误差 (empirical error) 是 h 的预测值 (prediction) 与 X 中给定的预期值 (required value) 不同的训练实例所占的比例。对于给定的训练集 X , 假设 h 的误差是

$$E(h|X) = \sum_{i=1}^N 1(h(x^i) \neq r^i) \quad (2.6)$$

其中, 当 $a \neq b$ 时 $1(a \neq b)$ 为 1, 当 $a = b$ 时 $1(a \neq b)$ 为 0 (参见图 2-3)。

在我们的例子中, 假设类 \mathcal{H} 是所有可能的矩形的集合。每个四元组 $(p_1^h, p_2^h, e_1^h, e_2^h)$ 都定义 \mathcal{H} 中的一个假设 h , 而我们需要选择其中最好的一个; 换句话说, 给定训练集, 我们需要找出这四个参数的值, 使得它涵盖所有的正例而不包括任何的负例。注意, 如果 x_1 和 x_2 是实数, 则存在无穷多个 h 满足上述条件, 也就是说, 对于这些 h 误差 E 为零。但是, 给定一个接近于正例和负例边界的某个未来实例, 不同的候选假设可能做出

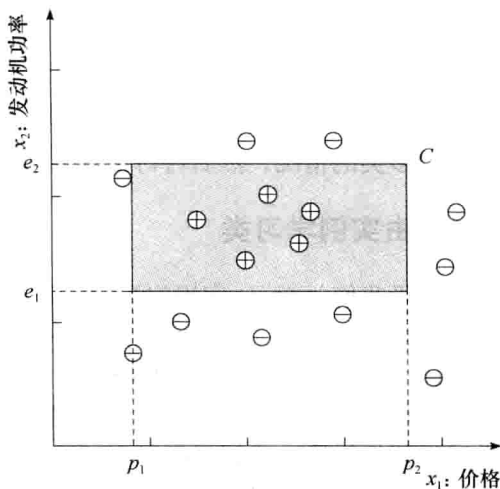


图 2-2 假设类的实例。家用汽车类是价格和发动机功率空间中的矩形

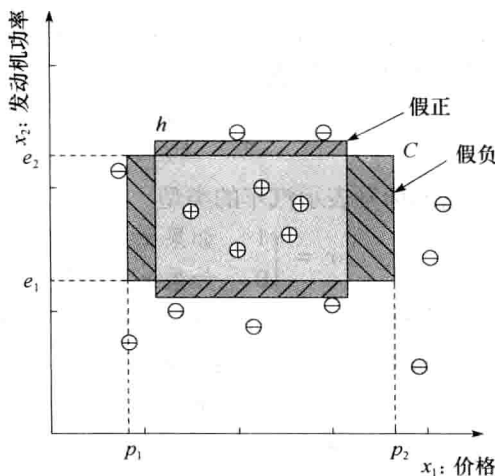


图 2-3 C 是实际的类, h 是我们的诱导假设。 C 为 1 而 h 为 0 的点为假负, C 为 0 而 h 为 1 的点为假正。其他点, 即真正和真负, 都被正确地分类

不同的预测。这是泛化 (generalization) 问题, 即我们的假设对不在训练集中的未来实例的分类的准确率如何。

一种可能的策略是找出最特殊的假设 (most specific hypothesis) S , 它是涵盖所有正例而不包括任何负例的最紧凑的矩形 (参见图 2-4)。这样就得出一个假设 $h = S$, 作为我们的诱导类 (induced class)。注意实际的类 C 可能会比 S 更大但绝不会更小。最一般的假设 (most general hypothesis) G 是涵盖所有正例而不包括任何负例的最大矩形 (参见图 2-4)。对于任何介于 S 和 G 之间的 $h \in \mathcal{H}$, h 为无误差的有效假设, 称作与训练集相容 (consistent), 并且这样的 h 形成解空间 (version space) 给定另一个训练集, S 、 G 、解空间、参数, 因此学习得到的假设 h 可能不同。

实际上, 依赖于训练集和假设类, 可能存在多个 S_i 和 G_i , 它们分别形成 S -集和 G -集。 S -集中的每个假设都与所有的实例相容, 并且不存在更特殊的相容假设。类似地, G -集中的每个假设都与所有的实例相容, 并且不存在更一般的相容假设。这两个集合形成边界集, 而它们之间的任何假设都是相容的, 并且是解空间的一部分。存在一个称作候选删除的算法, 随着逐个看到训练实例, 它增量地更新 S -集和 G -集, 参见 Mitchell 1997。我们假定 \mathcal{X} 足够大, 则存在唯一的 S 和 G 。

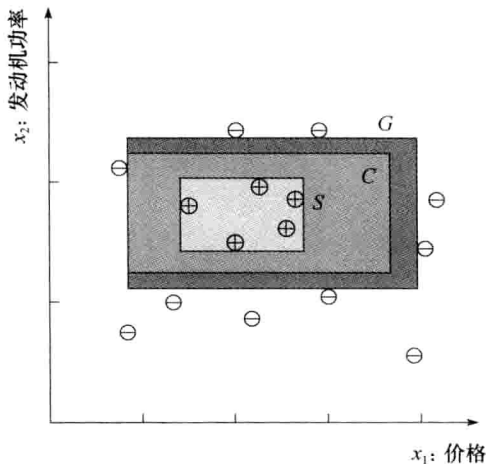
给定 \mathcal{X} , 我们可以找到 S 或 G , 或解空间中的任意 h , 并将它作为我们的假设 h 。直观地, h 应该选取 S 与 G 的中间, 这将增大边缘 (margin), 而边缘是边界和与它最近的实例之间的距离 (参见图 2-5)。为了使我们的误差函数在具有最大边缘的 h 上最小化, 我们应该选择这样的误差 (损失) 函数, 它不仅检查实例是否在边界的正确一侧, 而且还要指出实例离边界多远。也就是说, 取代返回 0 或 1 的 $h(x)$, 我们需要一个返回携带 x 到边界距离度量的假设, 并且需要一个使用该值、不同于检查相等性 $1(\cdot)$ 的损失函数。

在某些应用中, 错误决策的代价可能很高, 并且任何 S 和 G 之间的实例都是不确定的 (doubt) 实例, 由于缺乏数据支持, 这些不确定实例无法被确定地标注。在这种情况下, 系统将拒绝考虑 (reject) 这些实例, 并留待人类专家判定。

这里, 我们假定 \mathcal{H} 包含 C , 即存在 $h \in \mathcal{H}$, 使得 $E(h | \mathcal{X})$ 为 0。给定假设类 \mathcal{H} , 可能存在不能学习 C 的情况, 即不存在 $h \in \mathcal{H}$, 使得误差为 0。因此, 对于任何的应用, 我们都需要确信 \mathcal{H} 有足够的柔性, 或 \mathcal{H} 具有足够的“能力”学习 C 。

2.2 VC 维

假定我们有一个数据集, 包含 N 个点。这 N 个点可以用 2^N 种方法标记为正例和负例。因此, N 个数据点可以定义 2^N 种不同的学习问题。如果对于这些问题中的任何一个, 我们都能够找到一个假设 $h \in \mathcal{H}$ 将正例和负例分开, 那么我们就称 \mathcal{H} 散列 (shatter) N 个点。也就是说, 可以用 N 个点定义的任意的学习问题都能够用一个从 \mathcal{H} 中抽取的假设无误差地学习。



24

25

26

可以被 \mathcal{H} 散列的点的最大数量称为 \mathcal{H} 的 VC 维 (Vapnik-Chervonenkis dimension), 记为 $VC(\mathcal{H})$, 它度量假设类 \mathcal{H} 的学习能力 (capacity)。

在图 2-6 中, 我们可以看到, 轴平行的矩形能够散列二维空间的 4 个点。因此, 当 \mathcal{H} 为二维空间中轴平行的矩形的假设类时, $VC(\mathcal{H})$ 等于 4。在计算 VC 维时, 能找到 4 个被散列的点就够了; 没有必要去散列二维空间中任意 4 个点。例如, 位于同一直线上的 4 个点不能被矩形散列。然而, 我们无法在二维空间的任何位置设置 5 个点, 使得对于所有可能的标记, 一个矩形能够分开正例和负例。

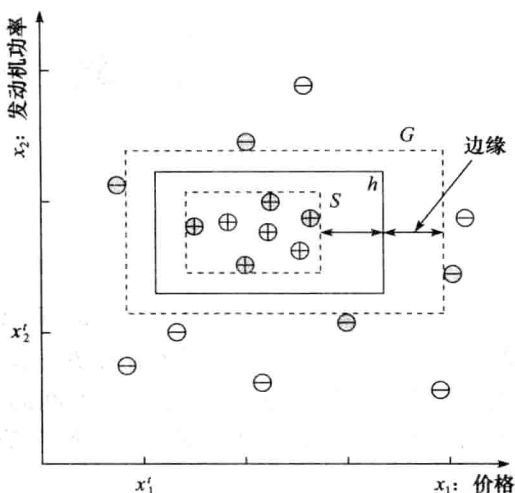


图 2-5 为了获得最佳分离, 我们选择具有最大边缘的假设。带阴影的实例是定义 (或支撑) 边缘的实例; 可以删除其他实例, 而不会影响 h

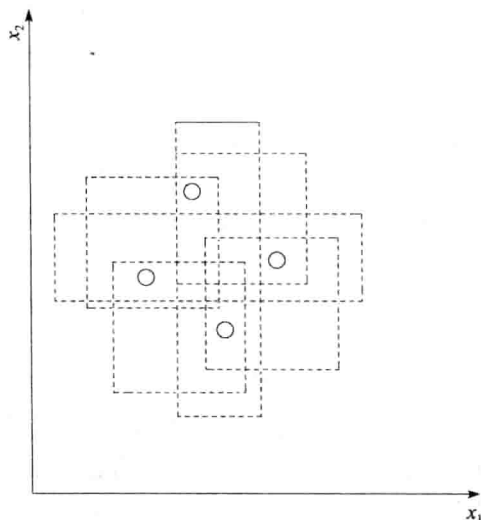


图 2-6 轴平行的矩形能够散列 4 个点, 其中只显示了覆盖两个点的矩形

也许 VC 维看起来比较悲观, 它告诉我们使用矩形作为假设类, 我们只能学习包括 4 个点的数据集。能够学习含有 4 个点的数据集的学习算法不是很有用。然而, 这是因为 VC 维独立于数据实例的概率分布。在实际生活中, 世界是平滑变化的, 在大多数时间相近的实例具有相同的标记, 我们并不需要担心所有可能的标记。有很多包含远不止 4 个点的数据集都可以通过我们的假设类来学习 (参见图 2-1)。因此, 即便是具有较小 VC 维的假设类也是有应用价值的, 并且比那些较大的 VC 维 (例如, 具有无穷 VC 维的查找表) 更可取。

2.3 概率逼近正确学习

使用最紧凑的矩形 S 作为假设, 希望找出我们需要多少实例。我们希望我们的假设是近似正确的, 即误差概率不超过某个值。还要对我们的假设有信心, 因为我们想知道我们的假设在大多数时间里都是正确的。因此我们希望我们的假设 (以我们可以指定的概率) 是正确的。

在概率逼近正确 (Probably Approximately Correct, PAC) 学习中, 给定类 C 和从未知但具

有确定概率分布 $p(x)$ 中抽取的样本, 我们希望找出样本数 N , 使得对于任意的 $\delta \leq 1/2$ 和 $\varepsilon > 0$, 假设 h 的误差至多为 ε 的概率至少为 $1 - \delta$ 。

$$P\{C\Delta h \leq \varepsilon\} \geq 1 - \delta$$

其中, $C\Delta h$ 是 C 与 h 不同的区域。

在这种情况下, 因为 S 是最紧凑的可能的矩形, C 与 $h=S$ 之间的误差区域是四个矩形条带之和(参见图 2-7)。我们希望确保正例落在该区域(导致错误)的概率最多为 ε 。对于任何这样的条带, 如果我们能够确保其概率上界为 $\varepsilon/4$, 则误差最多为 $4(\varepsilon/4) = \varepsilon$ 。注意, 我们将矩形角部的重叠部分计算了两次, 并且这种情况下总的实际误差小于 $4(\varepsilon/4)$ 。随机抽取的样本不在此条带中的概率是 $1 - \varepsilon/4$ 。所有 N 个独立抽取的样本不在此条带中的概率为 $(1 - \varepsilon/4)^N$, 所有 N 个独立抽取的样本不在任意这四个矩形条带中的概率最多为 $4(1 - \varepsilon/4)^N$, 我们希望其最大值为 δ 。我们有不等式

$$(1 - x) \leq \exp[-x]$$

因此, 如果选定 N 和 δ 满足

$$4\exp[-\varepsilon N/4] \leq \delta$$

则我们有 $4(1 - \varepsilon/4)^N \leq \delta$ 。不等式两边同时除以 4, 再取(自然)对数, 并重新排列各项, 我们得到

$$N \geq (4/\varepsilon) \log(4/\delta) \quad (2.7)$$

因此, 只要我们至少从 C 中取 $(4/\varepsilon) \log(4/\delta)$ 个独立样本, 并使用紧凑矩形作为我们的假设 h , 则在置信概率(confidence probability)至少为 $1 - \delta$ 的情况下, 一个给定点被误分类的错误概率最多为 ε 。减少 δ 我们可以有任意大的置信度, 而减少 ε 我们可以有任意小的误差, 并且我们在不等式(2.7)中看到, 样本的数量是分别随 $1/\varepsilon$ 和 $1/\delta$ 呈线性和对数缓慢增长的函数。

2.4 噪声

噪声(noise)是数据中有害的异常。由于噪声的存在, 类的学习可能更加困难, 并且使用简单的假设可能做不到零误差(参见图 2-8)。噪声有以下几种解释:

- 记录输入属性可能不准确, 这可能导致数据点在输入空间的移动。

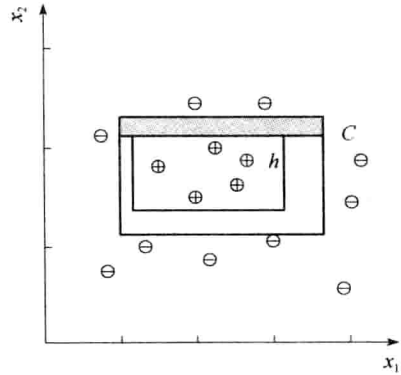


图 2-7 h 与 C 之差是四个矩形条带之和, 其中一个用阴影显示

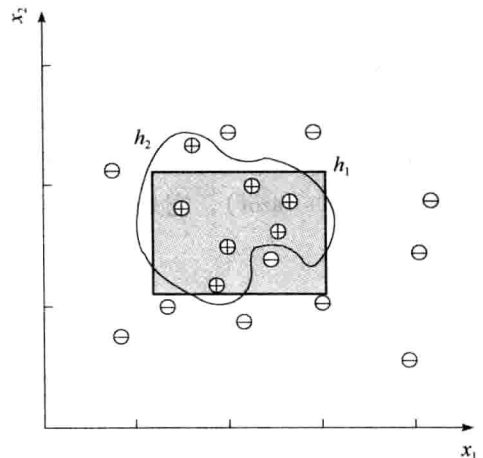


图 2-8 当有噪声时, 在正例和负例之间不存在一个简单的边界, 使用简单假设也许不可能达到零误差的分类结果。矩形是具有 4 个定义隅角的参数的简单假设。使用大量控制点的分段函数能够导出任意的封闭图形

- 标记数据点可能有错，可能将正例标记为负的，或相反。这种情况称为指导噪声 (teacher noise)。
- 可能存在我们没有考虑到的附加属性，而它们会影响实例的标注。这些附加属性可能是隐藏的 (hidden) 或潜在的 (latent)，因此是不可观测的。这些被忽略的属性所造成的影响作为随机成分，是“噪声”的一部分。

如图 2-8 所示，当有噪声时，在正负实例之间不存在简单的边界，并且为了将它们分开，我们需要对应于具有更大能力的假设类的复杂假设。矩形可以用 4 个数定义，然而为了定义更复杂的形状，我们就需要具有大量参数的更复杂的模型。利用这些复杂模型，我们可以更好地拟合数据，得到零误差 (参见图 2-8 中的曲线图形)。另一个可行的方法是保持模型的简单性并允许一些误差的存在 (参见图 2-8 中的矩形)。

使用简单的矩形 (除非其训练误差很大) 更有意义，原因如下：

1) 矩形是一种容易使用的简单模型。容易检查一个点是在矩形内还是在矩形外，并且对于未来的数据实例，我们都可以容易地检查它是正例还是负例。

2) 矩形是一种容易训练的简单的模型，并且具有较少的参数。相对任意图形的控制点来说，比较容易找到矩形的隅角值。利用小规模训练集，当训练实例有少许差异时，我们预期简单模型比复杂模型变化小一些：简单模型具有更小的方差 (variance)。另一方面，太简单的模型假设更多、更严格，并且如果潜在类并非那么简单，模型预测就可能失败：较简单的模型具有较大的偏倚 (bias)。求解最优模型相当于最小化偏倚和方差。

3) 矩形是容易解释的简单模型。矩形简单地对应在两个属性上定义的区间。通过学习简单的模型，我们能够从给定训练集的原始数据中提取信息。

4) 如果输入数据中确实存在错误标记的实例或噪声，并且实际的类确实就是像矩形这样的简单模型，那么由于矩形具有较小的方差，并且较少地被单个实例所影响，所以尽管简单矩形可能导致训练集上较大的误差，它也是比曲线图形更好的分类器。我们说简单 (但不是太简单的) 模型会比复杂模型泛化能力更好。该规则就是著名的奥克姆剃刀规则 (Occam's razor)，它是说较简单的解释看上去更可信，并且任何不必要的复杂性都应该被摒弃。

2.5 学习多类

在前面的家用汽车例子中，我们有属于家用汽车类的正例和属于其他所有汽车类别的负例。这是一个两类 (two-class) 问题。通常情况下，我们有 K 个类，记为 $C_i, i = 1, \dots, K$ ，并且每个输入实例严格地属于其中一个类。训练集形如

$$X = \{\mathbf{x}^t, \mathbf{r}^t\}_{t=1}^N$$

其中 \mathbf{r} 是 K 维的，并且

$$r_i^t = \begin{cases} 1 & \text{如果 } \mathbf{x}^t \in C_i \\ 0 & \text{如果 } \mathbf{x}^t \in C_j, j \neq i \end{cases} \quad (2.8)$$

一个例子在图 2-9 中给出，其中实例来自三个类：家用汽车、运动汽车和豪华轿车。

在用于分类的机器学习中，我们希望学习将一个类与所有其他类分开的边界。这样，我

们把 K -类的分类问题看作是 K 个两类问题。属于 C_i 类的训练实例是假设 h_i 的正例，属于所有其他类的训练实例是假设 h_i 的负例。因此，在 K -类的分类问题中，我们要学习 K 个假设，使得

$$h_i(\mathbf{x}^t) = \begin{cases} 1 & \text{如果 } \mathbf{x}^t \in C_i \\ 0 & \text{如果 } \mathbf{x}^t \in C_j, j \neq i \end{cases} \quad (2.9)$$

整体经验误差对所有类在所有实例上的预测取和：

$$E(\{h_i\}_{i=1}^K | \mathcal{X}) = \sum_{t=1}^N \sum_{i=1}^K 1(h_i(\mathbf{x}^t) \neq r_i^t) \quad (2.10)$$

在理想情况下，对于给定的 \mathbf{x} ，只有其中一个假设 $h_i(\mathbf{x}) (i=1, \dots, K)$ 为 1，并且我们能够选定一个类。但是，当没有或者有两个或更多的 $h_i(\mathbf{x})$ 为 1 时，我们就无法选定一个类，这是不确定的 (doubt) 情况并且分类器要拒绝这种情况。

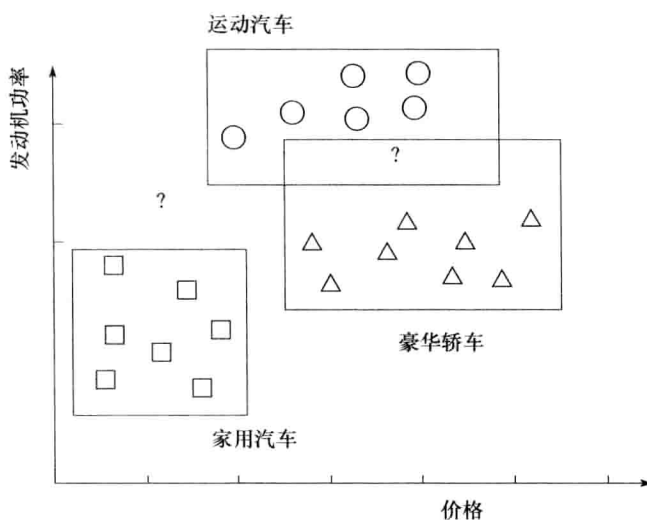


图 2-9 有三个类：家用汽车、运动汽车和豪华轿车。有三个归纳的假设，每个假设覆盖一个类的实例而不包括另外两个类的实例。“?”为拒绝区域，其中没有类或多个类被选中

在学习家用汽车的例子中，我们只用了一个假设，并且只对正例样本建模。任何未包括在其中的实例都不是家用汽车。作为另一种选择，有时我们可能更倾向于构建两个假设，一个是对正例，另一个是对负例。这也为被另一个假设所覆盖的负例假定一个结构。将家用汽车与运动汽车分开就是一个这样的问题，每个类都有其自己的结构。这种处理的优点在于，如果输入的是一辆豪华轿车，我们就能够通过两个假设来判定其为负例并丢弃该输入。

如果我们预料数据集中所有类的结构 (输入空间中的形状) 都类似，则可以对所有的类使用相同的假设。例如，在手写数字识别数据集中，我们预料所有数字都具有类似的分布。但是，在医疗诊断数据集中，例如，我们有病人和健康人两类，那么，这两类可能具有完全不同的分布。一个人是病人可能有不同原因，反映在输入中的不同：所有健康的人都是类似的，而每个病人都有他们自己的病情。

2.6 回归

在分类问题中，给定一个输入，所产生的输出是一个布尔值，这是一个是/否型答案。

34 当输出是数值型的值时, 我们所希望学习的不是一个类 $C(\mathbf{x}) \in \{0, 1\}$, 而是一个连续函数。在机器学习中, 函数是未知的, 不过我们有从其中抽取的训练集

$$\mathcal{X} = \{\mathbf{x}^i, r^i\}_{i=1}^N$$

其中 $r^i \in \mathfrak{R}$ 。如果不存在噪声, 则该任务是插值(interpolation)。我们希望找到通过这些点的函数 $f(\mathbf{x})$, 使得

$$r^i = f(\mathbf{x}^i)$$

在多项式插值(polynomial interpolation)中, 给定 N 个点, 我们找出可以用来预测任何 \mathbf{x} 的输出的 $(N-1)$ 阶多项式。如果 \mathbf{x} 落在训练集中 \mathbf{x}^i 的值域之外, 则该方法称为外插或外推(extrapolation)。例如, 在时间序列预测中, 我们拥有截至目前的数据, 而且希望预测未来的值。在回归(regression)分析中, 噪声添加到未知函数的输出上

$$r^i = f(\mathbf{x}^i) + \varepsilon \quad (2.11)$$

其中 $f(\mathbf{x}) \in \mathfrak{R}$ 是未知函数, ε 是随机噪声。关于噪声的解释是, 存在我们无法观察到的额外的隐藏(hidden)变量

$$r^i = f^*(\mathbf{x}^i, \mathbf{z}^i) \quad (2.12)$$

其中 \mathbf{z}^i 表示这些隐藏变量。我们希望通过模型 $g(\mathbf{x})$ 来逼近输出。训练集 \mathcal{X} 上的经验误差是

$$E(g|\mathcal{X}) = \frac{1}{N} \sum_{i=1}^N [r^i - g(\mathbf{x}^i)]^2 \quad (2.13)$$

因为 r 和 $g(\mathbf{x})$ 是数值量(例如, 属于 \mathfrak{R}), 因此存在定义在其值域上的序, 并且我们可以定义值之间的距离(distance)为差的平方。相对于分类使用的等于或不等于来说, 距离给我们提供了更多的信息。差的平方是一种可以使用的误差函数, 另一种误差函数是差的绝对值。在后续章节中, 我们将会看到一些其他例子。

我们的目标是找到最小化经验误差的 $g(\cdot)$ 。我们的方法又是相同的, 我们对 $g(\cdot)$ 假定一个的具有少量参数的假设类。如果假定 $g(\mathbf{x})$ 是线性的, 则我们有

$$35 \quad g(\mathbf{x}) = w_1 x_1 + \cdots + w_d x_d + w_0 = \sum_{j=1}^d w_j x_j + w_0 \quad (2.14)$$

现在, 再回到 1.2.3 节的例子, 在那里我们估计一辆二手车的价格。当时我们使用单个输入的线性模型

$$g(x) = w_1 x + w_0 \quad (2.15)$$

其中, w_1 和 w_0 是需要从数据中学习的参数。 w_1 和 w_0 的值应该使下式最小化

$$E(w_1, w_0|\mathcal{X}) = \frac{1}{N} \sum_{i=1}^N [r^i - (w_1 x^i + w_0)]^2 \quad (2.16)$$

其最小点可以通过求 E 关于 w_1 和 w_0 的偏导数, 令偏导数为 0, 并求解这两个未知量来计算:

$$w_1 = \frac{\sum_i x^i r^i - \bar{x} \bar{r} N}{\sum_i (x^i)^2 - N \bar{x}^2} \quad (2.17)$$

36

$$w_0 = \bar{r} - w_1 \bar{x}$$

其中, $\bar{x} = \sum_i x^i / N$, $\bar{r} = \sum_i r^i / N$ 。找到的直线如图 1-2 所示。

如果线性模型过于简单，它就会太受限制，导致大的近似误差，并且在这种情况下，输出可以取输入的较高阶的函数，例如二次函数

$$g(x) = w_2x^2 + w_1x + w_0 \quad (2.18)$$

其中类似地，我们有参数的解析解。当多项式的阶增加时，训练数据上的误差将会降低。但是高阶多项式关注个体样本，而不是捕获数据一般趋势(参看图 2-10 中的六次多项式)。因此，当精确调整的模型复杂性达到潜在数据的函数的复杂度时，我们应该谨慎行事。

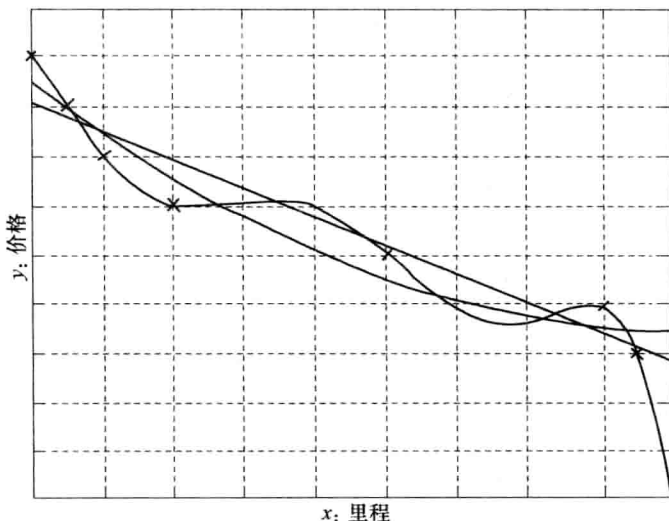


图 2-10 拟合相同的数据点集的线性、二次和六次多项式。最高阶的多项式(六次多项式)给出了正确的拟合，但是给定更多数据，真实的曲线很可能不是这种形状。二次多项式看起来比线性拟合好，它捕获了训练数据的走势

2.7 模型选择与泛化

我们用从实例学习布尔函数作为例子开始。在布尔函数中，所有的输入和输出均为二元的。 d 个二元值有 2^d 种可能的写法。因此，对于 d 个输入，训练集最多能有 2^d 个样本实例。如表 2-1 所示，其中的每一位都能标记为 0 或 1，因而对于 d 个输入，将有 2^{2^d} 个可能的布尔函数。

表 2-1 2 个输入存在 4 种可能的情况和 16 种可能的布尔函数

x_1	x_2	h_1	h_2	h_3	h_4	h_5	h_6	h_7	h_8	h_9	h_{10}	h_{11}	h_{12}	h_{13}	h_{14}	h_{15}	h_{16}
0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1
0	1	0	0	0	0	1	1	1	1	0	0	0	0	1	1	1	1
1	0	0	0	1	1	0	0	1	1	0	0	1	1	0	0	1	1
1	1	0	1	0	1	0	1	0	1	0	1	0	1	0	1	0	1

每一个不同的训练样本都会去掉一半的假设，即去掉那些猜测出错的假设。例如，假定我们有 $x_1 = 0$, $x_2 = 1$ ，而输出为 0，这种情况就去掉了假设 h_5 、 h_6 、 h_7 、 h_8 、 h_{13} 、 h_{14} 、 h_{15} 和 h_{16} 。这是观察学习的一种途径；随着我们看到更多的训练样例，我们逐步去掉那些与训练

数据不一致的假设。在布尔函数的情况下，为了最终得到单个假设，我们需要看到所有的 2^d 个训练样本。如果给定的训练集只包含所有可能实例的一个小子集（通常情况就是如此），也就是说，如果我们仅对少量情况知道输出应该是什么，则解是不唯一的。看到 N 个样本实例后，还有 2^{2^d-N} 个可能的函数。这是一个不适定问题（ill-posed problem），其中仅靠数据本身不足以找到唯一解。

在其他的学习应用中，在分类、回归中也存在同样的问题。随着我们看到更多的训练实例，我们对潜在函数的了解就更多，并且我们从假设类去掉更多不一致的假设，但是我们还剩下许多一致的假设。

这样，由于学习是一个不适定问题，并且单靠数据本身不足以找到解，因此我们应该做一些特别的假设，以便得到已有数据的唯一解。我们把为了使学习成为可能所做的假设集称为学习算法的归纳偏倚（inductive bias）。引入归纳偏倚的一种途径是假定一个假设类。在学习家用汽车类时，存在着无限种将正例与负例分开的方法。假定矩形是一种归纳偏倚，而后假定最紧凑的矩形就是另外一种归纳偏倚。在线性回归中，假定线性函数也是一种归纳偏倚，而在所有直线中选择最小化平方误差的直线则是另一种归纳偏倚。

然而，我们知道，每个假设类都有一定的能力，并且只能够学习确定的函数。使用具有更大能力、包含更复杂的假设的假设类，可以扩充可学习的函数类。例如，假设类 \mathcal{H} “两个互不重叠的矩形的并”具有更大的能力，但是其假设也更加复杂。类似地，在回归分析中，随着多项式的阶增高，其能力和复杂性也不断增加。现在的问题是要决定在哪里停止。

因此，如果没有归纳偏倚，则学习将是不可能的，而且现在的问题是如何选择正确的偏倚。该问题称作模型选择（model selection），即在可能的模型 \mathcal{H} 之间选择。对于这种问题的解答，我们应当记住机器学习的目标很少是去复制训练数据，而是预测新情况。也就是说，我们希望对于训练集之外的输入（其正确的输出并没有在训练集中给出）能够产生正确的输出。训练集上训练的模型在多大程度上能够对新的实例预测出正确输出称为泛化（generalization）。

对最好的泛化来说，我们应当使假设的复杂性与潜在数据的函数的复杂性相匹配。如果 \mathcal{H} 没有函数复杂，例如，当试图用直线拟合从三次多项式抽取的数据时，则是欠拟合（underfitting）。这种情况下，随着复杂性的增高，训练误差和确认误差都会降低。但是，如果 \mathcal{H} 太过复杂，数据不足以约束该假设，我们最后也许得到不好的假设 $h \in \mathcal{H}$ 。例如，当用两个矩形拟合从一个矩形抽取的数据时，这种情况就会发生。或者如果存在噪声，则过分复杂的假设可能不仅学习潜在的函数，而且也学习数据中的噪声，导致很差的拟合。例如，用六次多项式拟合从三次多项式抽样的噪声数据时，这种情况就会发生。这称为过（分）拟合（overfitting）。在这种情况下，拥有更多的训练数据是有帮助的，但是只能达到一定程度。给定训练集和 \mathcal{H} ，我们可以找到最小化训练误差的 $h \in \mathcal{H}$ 。但是，如果 \mathcal{H} 选择不好，则无论选择什么样的 $h \in \mathcal{H}$ 都得不到好的泛化。

我们可以引用三元权衡（triple trade-off）（Dietterich 2003）来总结我们的讨论。在所有的由实例数据训练的学习算法中，存在以下三种因素之间的平衡：

- 拟合数据的假设的复杂性，即假设类的能力；
- 训练数据的总量；
- 在新的实例上的泛化误差。

随着训练数据量的增加，泛化误差降低。随着模型类 \mathcal{H} 复杂性的增加，泛化误差先降低，然后开始增加。过于复杂的 \mathcal{H} 的泛化误差可以通过增加训练数据的总量来控制，但是只能达到一定程度。如果数据由直线采样并用高阶多项式拟合，则在周围有训练数据的地方，拟合将被束缚在该直线附近；而在没有被训练的地方，高阶多项式的行为可能难以预测。

如果我们访问训练集以外的数据，我们就能够度量假设的泛化能力，即归纳偏倚的质量。我们通过将训练集划分为两个部分来模拟这一过程。我们使用一部分作训练（即找出假设），剩下的部分称作验证集（validation set），并用来检验假设的泛化能力。也就是说，给定可能的假设类集合 \mathcal{H}_i ，则我们可以在训练集上拟合最佳的 $h_i \in \mathcal{H}_i$ 。假定训练集和验证集都足够大，则在验证集上最精确的假设就是最好的假设（即具有最佳归纳偏倚的假设）。该过程称为交叉验证（cross-validation）。例如，为了找出多项式回归的正确的阶，给定多个不同阶的候选多项式，其中不同阶的多项式对应于不同的 \mathcal{H}_i ，我们在训练集上求出它们的系数，在验证集上计算它们的误差，并取具有最小验证误差的多项式作为最佳多项式。

注意，如果需要报告反映我们的最佳模型的期望误差的话，我们就不应该使用验证误差。我们已经使用验证集来选择最佳模型，并且它实际上已经成为训练集的一部分。我们需要第三个数据集——检验集（test set），有时也称为发布集（publication set），它包含在训练或验证阶段未使用过的数据。现实生活也有类似的情况，例如我们选修一门课程：老师在讲授一门课时，课堂上求解的例题构成了训练集，考试题目就是验证集，而我们在职业生涯中解决的问题则是检验集。

我们也不能一直使用相同的训练或验证集划分，因为一旦使用一次，验证集实际上就成为了训练数据的一部分。这就像老师每年都使用相同的考试题一样；精明的学生会意识到不必听课，仅仅记住这些问题的答案即可。

一定要记住，我们使用的训练数据是一个随机样本。也就是说，对于相同的应用，如果我们多次收集数据，则将得到稍微不同的数据集，拟合的 h 也稍微不同，并且具有稍微不同的验证误差。或者，如果我们把固定的数据集划分成训练、验证和检验集，则根据划分方式的不同，我们会有不同的误差。这些微小的区别使得我们可以估计多大的差别可以看作显著的（significant）而非偶然的。也就是说，在假设类 \mathcal{H}_i 和 \mathcal{H}_j 之间进行选择时，我们将在大量训练集和验证集上多次使用它们，并且检查 h_i 与 h_j 的平均误差之差是否大于多个 h_i 之间的平均差。在第19章，我们将讨论如何设计机器学习实验，以利用有限的数据来回答我们的问题（例如，“最好的假设类是哪个？”），以及如何分析实验结果，使得我们可以得到最少受随机性影响的、统计显著的结论。

40
~
41

2.8 监督机器学习算法的维

现在，让我们来总结并归纳上述要点。我们有样本

$$\mathcal{X} = \{x^t, r^t\}_{t=1}^N \quad (2.19)$$

该样本是独立同分布的（independent and identically distributed, iid）；次序并不重要，而所有的实例都取自相同的联合分布 $p(x, r)$ 。 t 指示 N 个实例中的一个， x^t 是任意维的输入，而 r^t 是相关联的预期输出。对于两类学习， r^t 是0/1；对于 K （ $K > 2$ ）类问题， r^t 是一个 K 维二元向量（其中恰有一维为1，其他各维均为0）；在回归分析中， r^t 是一个实数值。

我们的目标是使用模型 $g(x' | \theta)$ 来构建一个 r' 的好的、有用的近似。为了达到预期目标,我们必须做出三个决定:

1) 学习所使用的模型(model), 记作

$$g(x | \theta)$$

其中, $g(\cdot)$ 是模型, x 是输入, θ 是参数。 $g(\cdot)$ 定义假设类, 而 θ 的特殊值示例了假设类中的一个假设。例如, 在类的学习中, 我们把矩形当作模型, 其 4 个坐标值构成了 θ 。在线性回归中, 模型是输入的线性函数, 其斜率和截距是从数据中学习的参数。模型(归纳偏倚)由机器学习系统的设计者根据其应用知识背景决定, 参数由学习算法, 利用取样于实际应用的训练集而进行调整。

41 2) 损失函数(loss function) $L(\cdot)$, 用于计算预期输出 r' 与给定参数 θ 的当前值时我们对它的近似 $g(x' | \theta)$ 之间的差。逼近误差(approximation error)或损失(loss)是各单个实例之上的损失之和

$$E(\theta | \chi) = \sum_i L(r'_i, g(x'_i | \theta)) \quad (2.20)$$

在输出为 0/1 类的学习中, $L(\cdot)$ 检测相等或不等; 在回归分析中, 由于输出是数值型值, 我们有关于距离的序信息, 而且一种可能性是使用差的平方。

3) 最优化过程(optimization procedure)求解最小化近似误差的 θ^*

$$\theta^* = \arg \min_{\theta} E(\theta | \chi) \quad (2.21)$$

其中 $\arg \min$ 返回使 E 最小化的参数值。在回归中, 我们能够解析地求解最优化问题。使用更复杂的模型和误差函数, 我们需要使用更加复杂的优化方法。例如, 基于梯度的方法、模拟退火或遗传算法等。

为了做好上述工作, 必须满足以下条件: 首先, $g(\cdot)$ 的假设类应当足够大, 即要有足够的容量, 以便包含在含噪声情况下产生 r' 表示的数据的未知函数。第二, 必须有足够的训练数据, 使得我们从假设类中识别正确(或足够好)的假设。第三, 给定训练数据, 我们应当有好的优化方法, 以便找出正确的假设。

不同的机器学习方法之间的区别或者在于它们假设的模型(假设类/归纳偏倚)不同, 或者在于它们所使用的损失度量不同, 或者在于它们所使用的最优化过程不同。我们将在后续的章节中看到更多的例子。

2.9 注释

Mitchell 提出了解空间和候选排除算法, 使得当样本实例依次给出时, 可以增量地构建 S 和 G , 近期的评述可参见 Mitchell 1997。矩形学习取自 Mitchell 1997 的习题 2.4。
42 Hirsh(1990)讨论了当实例样本受到少量噪声影响时, 如何处理解空间。

有关机器学习最早的研究工作之一是 Winston(1975)提出的“几乎错过”(near miss)思想。几乎错过是一个与正例非常相似的负例。用我们的术语, 几乎错过就是可能落在 S 与 G 之间灰色区域的实例, 该实例将会影响边缘, 因而相对于普通的正例和负例来说, 它们对学习可能更有用。靠近边界的实例是定义(或支撑)边界的实例; 添加或删除那些被许多具有相同标号包围的实例则不会影响边界。

与此相关的思想是主动学习 (active learning), 其中学习算法能够自己生成实例, 并要求被标记, 而不像在其他学习算法中那样被动地被给定 (Angluin 1988) (参见习题 4)。

VC 维在 20 世纪 70 年代初由 Vapnik 和 Chervonenkis 提出, 新近的相关资源是 Vapnik 1995, 其中他指出“没有什么比好的理论更实用”。像在其他科学领域一样, 这也在机器学习领域得到了证实。你不必急于使用计算机; 你可以使用纸张、铅笔, 也许还需要橡皮擦之类的东西来思考, 节省自己的时间, 避免无用的编程。

PAC 模型由 Valiant (1984) 提出, 对于学习矩形的 PAC 分析来自 Blumer 等 (1989)。一本涵盖 PAC 学习及 VC 维的计算学习理论的好教材是 Kearns 和 Vazirani (1994)。

2.10 习题

- 假定我们的假设类是圆而不是矩形。参数是什么? 这种情况下如何计算圆假设的参数? 如果是椭圆又如何? 为什么用椭圆代替圆会更有意义? 如何将算法推广到 $K > 2$ 个类情况?
- 设想我们的假设类不是一个矩形而是两个 (或 $m > 1$ 个) 矩形的并, 请问这种假设类优点何在? 说明使用足够大的 m , 任何类都能够由这种假设类表示。
- 大部分学习算法的复杂度都是训练集的函数。你能提出一个发现冗余实例的过滤算法吗?
- 如果我们拥有能够给任何实例 x 提供标记的指导者, 那么我们应当在哪里选择 x , 以便用较少的询问来进行学习?
- 在公式 (2.13) 中, 我们对实际值与估计值之差的平方求和。该误差函数是使用最频繁的误差函数, 但它只是可行的误差函数之一。由于它对差的平方求和, 所以它对于离群点不是鲁棒的。为了实现鲁棒回归 (robust regression), 更好的误差函数是什么?
- 请推导式 (2.17)。
- 假定我们的假设类是直线的集合, 并且我们利用直线来分开正例与负例, 而不是用矩形来界定正例, 并将负例留在矩形外 (参见图 2-11)。证明直线的 VC 维为 3。
- 证明在二维空间中, 三角形假设类的 VC 维为 7。
(提示: 为了最佳分割, 最好在某圆上设置 7 个等距离的点。)
- 假定像习题 7 那样, 假设类是直线的集合。写一个图 2-11 直线分割正例样本与负例样本误差函数, 它不仅最小化误分类数, 而且也最大化边缘。
- 噪声的一个来源是标号错误。你能提出一种方法, 找出很可能是误标记的数据点吗?

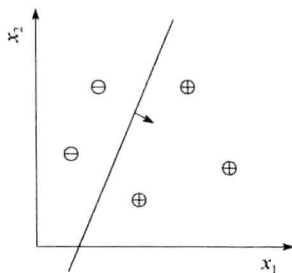


图 2-11 直线分割正例样本与负例样本

43

2.11 参考文献

- Angluin, D. 1988. "Queries and Concept Learning." *Machine Learning* 2: 319-342.
- Blumer, A., A. Ehrenfeucht, D. Haussler, and M. K. Warmuth. 1989. "Learnability and the Vapnik-Chervonenkis Dimension." *Journal of the ACM* 36: 929-965.

44

- Dietterich, T. G. 2003. "Machine Learning." In *Nature Encyclopedia of Cognitive Science*. London: Macmillan.
- Hirsh, H. 1990. *Incremental Version Space Merging: A General Framework for Concept Learning*. Boston: Kluwer.
- Kearns, M. J., and U. V. Vazirani. 1994. *An Introduction to Computational Learning Theory*. Cambridge, MA: MIT Press.
- Mitchell, T. 1997. *Machine Learning*. New York: McGraw-Hill.
- Valiant, L. 1984. "A Theory of the Learnable." *Communications of the ACM* 27: 1134-1142.
- Vapnik, V. N. 1995. *The Nature of Statistical Learning Theory*. New York: Springer.
- Winston, P. H. 1975. "Learning Structural Descriptions from Examples." In *The Psychology of Computer Vision*, ed. P. H. Winston, 157-209. New York: McGraw-Hill.

贝叶斯决策定理

我们讨论在不确定情况下决策的概率理论框架。在分类中，贝叶斯规则用来计算类的概率。我们将讨论推广到怎样做出合理的决策将期望风险最小化。我们还讨论从数据中学习关联规则。

3.1 引言

训练计算机使之根据数据进行推断是统计学和计算机科学的交叉领域，其中，统计学家提供由数据推断的数学框架，而计算机科学家研究推断方法在计算机上如何有效地实现。

数据来自一个不完全清楚的过程。将该过程作为随机过程建模表明我们缺乏知识。也许该过程是确定性的，但是因为我们没有获取关于它的完全知识的途径，所以我们把它作为一个随机过程建模，并且用概率理论来分析它。说到这，在继续阅读本章之前先跳到附录，温习基本的概率知识也许是一个好主意。

投硬币是一个随机的过程，因为我们不能够预测任意一次投币结果是正面还是反面（这就是为什么我们投币、买彩票或者买保险的原因）。我们只能谈论下一次投币是正面还是反面的概率。有证据显示，如果我们取得一些额外的数据，如硬币的确切成分，它的最初位置，投币的力量和投币的方向，何处以及如何接住等，则投币的准确结果就是可以预测的。

我们不能获取的那些额外的数据称为不可观测的变量(unobservable variable)。在投币的这个例子中，唯一可观测的变量(observable variable)是投币的结果。用 z 表示不可观测的变量， x 表示可观测的变量，事实上我们有

$$x = f(z)$$

其中， $f(\cdot)$ 是一个确定性函数，它定义不可观测数据的输出。因为我们不能用这种方式对该过程建模，所以我们定义输出 X 为指明该过程、由概率分布 $P(X=x)$ 抽取的随机变量。

投币的结果是正面或是反面，而我们定义一个随机变量，在两个值中取值。令 $X=1$ 代表投币的结果是正面， $X=0$ 代表投币结果是反面。 X 服从伯努利分布，其中参数 p_0 是投币结果为正面的概率。

$$P(X=1) = p_0 \text{ 并且 } P(X=0) = 1 - P(X=1) = 1 - p_0$$

假设要预测下一次投币的结果。如果我们知道 p_0 的值，则当 $p_0 > 0.5$ 时，预测将是正面，否则是反面。这是因为，如果选择更可能的情况，则错误的概率，即 1 减去选择的概率，将会最小。如果这是一个 $p_0 = 0.5$ 的公平投币，则我们没有比总是选择正面或者我们自己做公平投币更好的预测手段！

如果我们不知道 $P(X)$ ，并且想从给定的样本估计它，就需要统计学知识了。我们有一个样本 X ，包含由可观测变量 x^i 的概率分布(记作 $p(x)$)抽取出的样例。目的是使用样本 X

构造一个它的近似 $\hat{p}(x)$ 。

在投币例子中, 样本包含了 N 次投币的结果。然后利用 \mathbf{X} , 我们可以估计 p_0 。 p_0 是唯一定义该分布的参数。 p_0 的估计是

$$\hat{p}_0 = \frac{\#\{\text{结果为正面的投币}\}}{\#\{\text{投币}\}}$$

使用数值随机变量, 如果投币 t 的结果是正面, 则 x^t 为1, 否则为0。给定样本|正面, 正面, 反面, 正面, 反面, 反面, 正面, 正面|, 则 $\mathbf{X} = \{1, 1, 1, 0, 1, 0, 0, 1, 1\}$, 并且

$$\hat{p}_0 = \frac{\sum_{i=1}^N x^i}{N} = \frac{6}{9}$$

48

3.2 分类

在1.2.2节, 我们讨论了信用评分。在那里我们看到, 根据过去的交易, 银行的某些客户是低风险的, 因为他们还清了贷款并且银行从他们那里获利; 其他客户是高风险的, 因为他们不能偿还贷款。分析这些数据, 我们想学习“高风险客户”类, 使得未来有新的贷款申请时, 我们可以检查申请者是否符合“高风险客户”类的描述, 并据此决定接受还是拒绝该申请。使用关于申请的知识, 我们假定有两种信息是可观测的。我们观测它们是因为我们有理由相信它们为我们提供了客户信用信息。例如, 我们假定我们观测客户的年收入和存款, 它们分别用随机变量 X_1 和 X_2 表示。

可以断言, 如果我们能够获得客户的其他知识, 比如关于客户经济状况的全部细节和全部知识, 他/她的意图、道德规范等, 则我们可以确定地计算出客户是“低风险客户”还是“高风险客户”。但是, 这些是不可观测的, 而使用我们能够观测的信息, 客户的信用可以用观测条件 $\mathbf{X} = [X_1, X_2]^T$ 下的伯努利随机变量 C 表示, 其中 $C=1$ 表示高风险客户, $C=0$ 表示低风险客户。如果我们知道 $P(C | X_1, X_2)$, 则当一个 $X_1 = x_1$ 和 $X_2 = x_2$ 的新申请到达时, 我们可以

$$\text{选择} \begin{cases} C = 1 & \text{如果 } P(C = 1 | x_1, x_2) > 0.5 \\ C = 0 & \text{否则} \end{cases}$$

或等价地

$$\text{选择} \begin{cases} C = 1 & \text{如果 } P(C = 1 | x_1, x_2) > P(C = 0 | x_1, x_2) \\ C = 0 & \text{否则} \end{cases} \quad (3.1)$$

错误的概率是 $1 - \max(P(C = 1 | x_1, x_2), P(C = 0 | x_1, x_2))$ 。这个例子与抛硬币的例子类似, 只是伯努利随机变量 C 是在两个其他观测变量条件下的随机变量。我们用 \mathbf{x} 表示观测变量向量 $\mathbf{x} = [x_1, x_2]^T$ 。于是, 问题是要能够计算 $P(C | \mathbf{x})$ 。使用贝叶斯规则, 它可以表示为

49

$$p(C | \mathbf{x}) = \frac{P(C)p(\mathbf{x} | C)}{p(\mathbf{x})} \quad (3.2)$$

$P(C=1)$ 称为 C 取值1的先验概率(prior probability)。在我们的例子中, 它对应客户是

高风险客户的概率，而不管 \mathbf{x} 取什么值。称它为先验概率，是因为它是我们看到观测量 \mathbf{x} 之前就获得的关于 C 值的知识，满足

$$P(C = 0) + P(C = 1) = 1$$

$p(\mathbf{x} | C)$ 称为类似然 (class likelihood)，是属于 C 的事件具有相关联的观测值 \mathbf{x} 的条件概率。在我们的例子中， $p(x_1, x_2 | C = 1)$ 是高风险客户具有 $X_1 = x_1, X_2 = x_2$ 的概率。这就是通过数据我们得到的关于类的信息。

$p(\mathbf{x})$ 是证据 (evidence)，是看到观测 \mathbf{x} 的边缘概率，无论它是正实例还是负实例。

$$p(\mathbf{x}) = \sum_C p(\mathbf{x}, C) = p(\mathbf{x} | C = 1)P(C = 1) + p(\mathbf{x} | C = 0)P(C = 0) \quad (3.3)$$

使用贝叶斯规则，组合先验知识和数据告诉我们的，在看到观测 \mathbf{x} 之后，计算概念的后验概率 (posterior probability) $P(C | \mathbf{x})$ 。

$$\text{后验} = \frac{\text{先验} \times \text{似然值}}{\text{证据}}$$

由于用证据规范化，后验的和为 1：

$$P(C = 0 | \mathbf{x}) + P(C = 1 | \mathbf{x}) = 1$$

一旦得到后验概率，我们就可以使用式 (3.1) 进行决策。从现在起，我们假定我们知道先验和似然。在稍后的章节中，我们会讨论如何从训练样本估计 $P(C)$ 和 $p(\mathbf{x} | C)$ 。

在一般情况下，我们有 K 个互斥和穷举的类 $C_i, i = 1, \dots, K$ ；例如，在光学数字识别中，输入是一个位图图像，有 10 个类。我们得到先验概率满足：

$$P(C_i) \geq 0 \text{ 并且 } \sum_{i=1}^K P(C_i) = 1 \quad (3.4)$$

$p(\mathbf{x} | C_i)$ 是已知属于类 C_i ，看到 \mathbf{x} 作为输入的概率。类 C_i 的后验概率计算如下

$$P(C_i | \mathbf{x}) = \frac{p(\mathbf{x} | C_i)P(C_i)}{p(\mathbf{x})} = \frac{p(\mathbf{x} | C_i)P(C_i)}{\sum_{k=1}^K p(\mathbf{x} | C_k)P(C_k)} \quad (3.5) \quad \boxed{50}$$

而为了将误差最小化，贝叶斯分类器 (Bayes' classifier) 选择具有最高后验概率的类；即

$$\text{选择 } C_i \text{ 如果 } P(C_i | \mathbf{x}) = \max_k P(C_k | \mathbf{x}) \quad (3.6)$$

3.3 损失与风险

决策的好坏程度或代价可能不同。金融机构对一个贷款申请人做出决定时会把潜在的收益和损失考虑在内。接受一个低风险的申请人会增加收益，而拒绝一个高风险的申请人会减小损失。错误地接受一个高风险的申请人带来的损失与错误地拒绝一个低风险的申请人带来的潜在收益是不同的。这种情况在其他领域，如医疗诊断、地震预测等，显得更加至关重要并且是非常不对称的。

让我们定义动作 α_i 为把输入指派到类 C_i 的决策，而 λ_{ik} 为输入实际属于 C_k 时采取动作 α_i 导致的损失 (loss)。采取动作 α_i 带来的期望风险 (expected risk) 是

$$R(\alpha_i | \mathbf{x}) = \sum_{k=1}^K \lambda_{ik} P(C_k | \mathbf{x}) \quad (3.7)$$

并且我们选择具有最小风险的动作：

$$\text{选择 } \alpha_i \text{ 如果 } R(\alpha_i | \mathbf{x}) = \min_k R(\alpha_k | \mathbf{x}) \quad (3.8)$$

让我们定义 K 个动作 $\alpha_i, i=1, \dots, K$, 其中 α_i 是把 \mathbf{x} 指派到 C_i 的动作。在 0-1 损失 (zero-one loss) 这种特殊情况下, 其中

$$\lambda_{ik} = \begin{cases} 0 & \text{如果 } i = k \\ 1 & \text{如果 } i \neq k \end{cases} \quad (3.9)$$

所有正确的决策没有损失, 并且所有错误具有相同的代价。采取动作 α_i 的风险是

$$\begin{aligned} R(\alpha_i | \mathbf{x}) &= \sum_{k=1}^K \lambda_{ik} P(C_k | \mathbf{x}) \\ &= \sum_{k \neq i} P(C_k | \mathbf{x}) \\ &= 1 - P(C_i | \mathbf{x}) \end{aligned}$$

51

因为 $\sum_k P(C_k | \mathbf{x}) = 1$ 。所以为了将风险最小化, 我们选择最有可能的情况。在后面章节中, 为了简单起见, 我们一直假定这种情况, 并且选择具有最高后验的类, 但是注意, 这确实是一种特殊情况, 并且很少应用具有对称的 0-1 损失。在一般情况下, 由后验到风险并且采取将风险最小化的动作是一种简单的后处理。

在一些应用中, 错误的决策 (即误分类) 也许会有很高的代价。一般情况下, 如果自动系统对它的决策的把握较低, 则需要一个更复杂 (例如人工的) 决策。例如, 如果我们使用光学数字识别器来读取信封上的邮编号码, 则错误地识别邮政编码将导致信件被发送到一个错误的目的地。

在这种情况下, 我们定义一个附加的拒绝 (reject) 或疑惑 (doubt) 动作 α_{K+1} , 而 $\alpha_i (i=1, \dots, K)$ 是在类 $C_i (i=1, \dots, K)$ 上的通常的决策动作 (Duda、Hart 和 Stork 2001)。

一个可能的损失函数是

$$\lambda_{ik} = \begin{cases} 0 & \text{如果 } i = k \\ \lambda & \text{如果 } i = K+1 \\ 1 & \text{否则} \end{cases} \quad (3.10)$$

其中 $0 < \lambda < 1$ 是选择第 $(K+1)$ 个拒绝动作导致的损失。拒绝的风险是

$$R(\alpha_{K+1} | \mathbf{x}) = \sum_{k=1}^K \lambda P(C_k | \mathbf{x}) = \lambda \quad (3.11)$$

而选择类 C_i 的风险是

$$R(\alpha_i | \mathbf{x}) = \sum_{k \neq i} P(C_k | \mathbf{x}) = 1 - P(C_i | \mathbf{x}) \quad (3.12)$$

最优决策规则是

$$\begin{aligned} \text{选择 } C_i & \text{ 如果对于所有的 } k \neq i \text{ 有 } R(\alpha_i | \mathbf{x}) < R(\alpha_k | \mathbf{x}), \text{ 并且} \\ & R(\alpha_i | \mathbf{x}) < R(\alpha_{K+1} | \mathbf{x}) \end{aligned}$$

$$\text{拒绝} \quad \text{如果 } R(\alpha_{K+1} | \mathbf{x}) < R(\alpha_i | \mathbf{x}), i=1, \dots, K \quad (3.13)$$

给定式 (3.10) 的损失函数, 上式可以简化为

$$\begin{aligned} \text{选择 } C_i & \text{ 如果对于所有的 } k \neq i \text{ 有 } P(C_i | \mathbf{x}) > P(C_k | \mathbf{x}), \text{ 并且} \\ & P(C_i | \mathbf{x}) > 1 - \lambda \end{aligned}$$

52

$$\text{拒绝} \quad \text{否则} \quad (3.14)$$

当 $0 < \lambda < 1$ 时, 这个方法是有意义的: 当 $\lambda = 0$ 时, 我们总是拒绝; 拒绝和正确的分类是同样好的。当 $\lambda \geq 1$ 时, 我们从不拒绝; 拒绝与错误的代价相同甚至超过错误的代价。

3.4 判别式函数

分类也可以看作是实现一组判别式函数(discriminant function) $g_i(\mathbf{x}) (i=1, \dots, K)$ 使得我们

$$\text{选择 } C_i, \text{ 如果 } g_i(\mathbf{x}) = \max_k g_k(\mathbf{x}) \quad (3.15)$$

令

$$g_i(\mathbf{x}) = -R(\alpha_i | \mathbf{x})$$

我们可以重新给出贝叶斯分类器, 并且最大化判别式函数对应最小化条件风险。当我们使用 0-1 损失函数时, 我们有

$$g_i(\mathbf{x}) = P(C_i | \mathbf{x})$$

或者忽略公共规范化项 $p(\mathbf{x})$, 我们可以写为

$$g_i(\mathbf{x}) = p(\mathbf{x} | C_i) P(C_i)$$

这把特征空间划分成 K 个决策区域(decision region) $\mathcal{R}_1, \dots, \mathcal{R}_K$, 其中 $\mathcal{R}_i = \{\mathbf{x} | g_i(\mathbf{x}) = \max_k g_k(\mathbf{x})\}$ 。这些区域被决策边界(decision boundaries), 即特征空间中的曲面分隔开, 其中平局出现在最大判别式函数之间(参见图 3-1)。

当有两个类的时候, 我们可以定义单个判别式

$$g(\mathbf{x}) = g_1(\mathbf{x}) - g_2(\mathbf{x})$$

并且我们

$$\text{选择} \begin{cases} C_1 & \text{如果 } g(\mathbf{x}) > 0 \\ C_2 & \text{否则} \end{cases}$$

一个例子是两类学习问题, 其中正例可以表示为 C_1 , 负例表示为 C_2 。当 $K=2$ 时, 分类系统是一个两分器(dichotomizer), 当 $K \geq 3$ 时, 它是一个多分器(polychotomizer)。

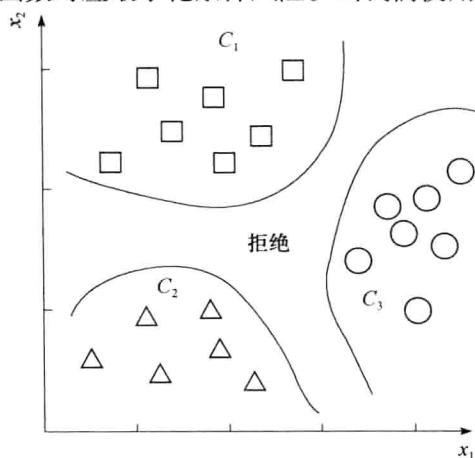


图 3-1 决策区域和决策边界的例子

3.5 效用理论

在式(3.7)中, 我们定义了期望风险并且选择最小化期望风险的动作。现在, 我们把它推广到效用理论, 它涉及我们对状态不确定时所做出的合理决策。我们假设给定证据 \mathbf{x} , 状态 S_k 的概率用 $P(S_k | \mathbf{x})$ 计算。我们定义一个效用函数(utility function) U_{ik} , 它度量当状态为 S_k 时采取动作 α_i 的好处。期望效用(expected utility)是

$$EU(\alpha_i | \mathbf{x}) = \sum_k U_{ik} P(S_k | \mathbf{x}) \quad (3.16)$$

一个合理的决策者选择最大化期望效用的动作

$$\text{选择 } \alpha_i, \text{ 如果 } EU(\alpha_i | \mathbf{x}) = \max_j EU(\alpha_j | \mathbf{x}) \quad (3.17)$$

对于分类, 决策对应选择类中的一个, 并且最大化期望效用等价于最小化期望风险。 U_{ik} 一般用货币术语度量, 并且这也为我们提供了一种定义损失矩阵 λ_{ik} 的途径。例如, 在我

54 们定义拒绝选项(3.10 式)时, 根据我们的特定应用, 如果我们知道一个正确决策能给我们带来多少钱的收益, 一个错误决策将使我们损失多少钱, 还有把决策委托给人类专家的代价多大, 那么我们就用货币单位, 而不是用 0、 λ 和 1 填写 U_{ik} 的正确值, 并且做出我们的决策来最大化期望收益。

对于拒绝, 我们在计算机程序做出的自动决策和代价较大但是正确可能性更高的人类决策之间选择。类似地, 可以设想一条多个自动决策者的流水线, 我们处理起来代价比较高, 但是正确的机会更大。我们将在第 17 章讨论这种级联以及如何组合多个学习器。

3.6 关联规则

关联规则 (association rule) 是形如 $X \rightarrow Y$ 的蕴涵式, 其中 X 是规则的前件 (antecedent), 而 Y 是规则的后件 (consequent)。关联规则的一个例子是购物篮分析, 通过购物篮分析, 我们希望发现项 X 和 Y 之间的依赖性。一个典型的应用是零售, 其中 X 和 Y 是销售的商品 (见 1.2.1 节)。

在学习关联规则时, 有三个度量需要计算:

- 关联规则 $X \rightarrow Y$ 的置信度 (confidence):

$$\text{Confidence}(X \rightarrow Y) \equiv P(Y|X) = \frac{P(X,Y)}{P(X)} = \frac{\#\{\text{购买 } X \text{ 和 } Y \text{ 的顾客}\}}{\#\{\text{购买 } X \text{ 的顾客}\}} \quad (3.18)$$

- 关联规则 $X \rightarrow Y$ 的提升度 (lift), 又称兴趣度 (interest)

55
$$\text{Lift}(X \rightarrow Y) = \frac{P(X,Y)}{P(X)P(Y)} = \frac{P(Y|X)}{P(Y)} \quad (3.19)$$

- 关联规则 $X \rightarrow Y$ 的支持度 (support):

$$\text{Support}(X,Y) \equiv P(X,Y) = \frac{\#\{\text{购买 } X \text{ 和 } Y \text{ 的顾客}\}}{\#\{\text{顾客}\}} \quad (3.20)$$

还有其他度量 (Omicinski 2003), 但是这三种, 特别是前两种被广泛认识和使用。置信度是我们通常计算的条件概率 $P(Y|X)$ 。为了能够说该规则具有足够的置信度, 它的值应该接近于 1, 并且显著大于人们购买 Y 的总体概率 $P(Y)$ 。我们也对最大化规则的支持度感兴趣, 因为即使有一个强置信度的依赖, 如果这样的顾客数量很少, 该规则也是没有价值的。支持度展示了规则的统计显著性, 而置信度展示了规则的强度。最小支持度和最小置信度值由公司设定, 并且从数据库中搜索具有更高支持度和置信度的所有规则。

如果 X 和 Y 是独立的, 则我们预料提升度接近于 1; 如果该比率的分子与分母不同, 即如果 $P(Y|X)$ 与 $P(Y)$ 不同, 则我们预料这两项之间存在依赖性; 如果提升度大于 1, 则我们可以说 X 使得 Y 更可能出现; 如果提升度小于 1, 则 X 使得 Y 更不可能出现。

这些公式可以很容易推广到多于两个项。例如, $\{X, Y, Z\}$ 是一个 3 项集, 而我们可以找出诸如 $X, Z \rightarrow Y$ 这样的规则, 即 $P(Y|X, Z)$ 。我们感兴趣的是找出具有足够高的支持度和置信度的所有规则, 并且由于销售数据库一般非常大, 所以我们希望通过少数几遍对于数据库的扫描找出它们。有一个有效的算法, 称作 Apriori 算法 (Agrawal 等, 1996), 来做这项工作。该算法分两步: (1) 找出频繁项集, 即找出具有足够支持度的项集; (2) 通过把频繁项集划分成两个子集, 分别作为前件和后件, 从而把频繁项集转换成具有足够置信度的规则。

1) 为了快速找出频繁项集(而不完全枚举项的所有可能子集), *Apriori* 算法利用如下事实: $\{X, Y, Z\}$ 要成为频繁的(具有足够支持度的), 它的所有子集 $\{X, Y\}$ 、 $\{X, Z\}$ 和 $\{Y, Z\}$ 也都应该是频繁的——添加另一项不会提高支持度。这就是说, 我们只需要检查其 2 项子集都频繁的 3 项集。换句话说, 如果知道一个 2 项集不是频繁的, 则可以剪掉它的所有超集而不必检查它们。

56

我们从找出频繁 1 项集开始, 并在每一步, 以归纳的方式由频繁 k 项集产生候选 $k+1$ 项集, 然后扫描数据, 以检查它们是否具有足够的支持度。为了方便访问, *Apriori* 把频繁项集存放在一个散列表中。注意, 随着 k 增加, 候选项集的个数将迅速减少。如果最长的频繁项集包含 n 项, 则总共需要扫描数据 $n+1$ 次。

2) 一旦我们找出了频繁 k 项集, 就需要把 k 项集划分成两个子集, 分别作为前件和后件, 把它转换成规则。正如产生频繁项集所做的那样, 我们从单项为后件、 $k-1$ 项为前件开始。因此, 对于所有可能的单项后件规则, 检查它是否具有足够的置信度, 如果没有就删除它。

注意, 对于相同的项集, 可能有多个以不同的子集为前、后件的规则。因此, 我们检查是否可以把一项从前件移到后件。后件中具有更多项的规则更特殊、更有用。这里, 与频繁项集产生一样, 我们利用如下事实: 为了得到一个后件中有两个项且具有足够置信度的规则, 则具有单项后件的两个规则本身都应该具有足够的置信度。也就是说, 从单项后件规则到两项后件规则不需要检查所有可能的两项后件(习题 7)。

应该记住, 规则 $X \rightarrow Y$ 不必蕴含因果关系, 而只是一种关联。在一个问题中, 可能还有一些隐藏变量, 它们的值不能通过证据获得。使用隐藏变量的优点是可以更容易定义依赖结构。例如, 在购物篮分析中, 当我们想找出销售商品之间的依赖时, 比如我们知道“婴儿食品”、“尿不湿”和“牛奶”之间的依赖性, 因为购买其中一种商品的顾客多半会买另两种。我们不表示这三者之间的依赖性, 而是指定一个隐藏结点“家有婴儿”作为这三种商品消费的隐藏原因, 我们将在第 16 章讨论的图模型使得我们可以表示这种隐藏变量。当存在隐藏结点时, 它们的值由给定的观测结点的值估计和填入。

57

3.7 注释

依据不确定数据进行决策已经有很长的历史, 并且人类一直在探索各种陌生领域, 寻找证据来排除不确定性: 例如天体、水晶球和咖啡杯。使用概率论, 从有意义的证据推理仅有几百年的历史。关于概率和统计学的历史, 以及拉普拉斯、伯努利和创建该理论的其他学者的一些早期论文见 Newman1988。

Russell 和 Norving(1995)给出了效用理论和信息值的很好讨论, 该书还用货币术语讨论了效用分配。Shafer 和 Pearl 1986 是不确定下推理的早期论文的汇集。

关联规则成功地用在许多数据挖掘应用中, 并且我们在许多网站上都看到这样的规则, 用来推荐书籍、电影、音乐等。算法很简单, 但是在大型数据库上的有效实现非常重要(Zhang 和 Zhang 2002; Li 2006)。稍后, 我们将会在第 16 章看到图模型如何把关联规则推广到非二元的情况, 其中关联可以有不同的类型, 也允许隐藏变量。

3.8 习题

1. 在两个类的问题中, 似然比(likelihood ratio)是

$$\frac{p(\mathbf{x} | C_1)}{p(\mathbf{x} | C_2)}$$

请用似然比写出判别函数。

2. 在两个类的问题中，对数几率(log odd)定义为

$$\log \frac{P(C_1 | \mathbf{x})}{P(C_2 | \mathbf{x})}$$

请用对数几率写出判别函数。

3. 在两类、两动作问题中，如果损失函数是 $\lambda_{11} = \lambda_{22} = 0$, $\lambda_{12} = 10$, $\lambda_{21} = 1$, 写出最佳决策规则。
4. 提出一种三级级联，当一级拒绝时，像式(3.10)那样使用下一级。如何在不同级中设定 λ ?
5. 某人做公平投币，如果结果是正面，你什么都得不到，否则就会得到 \$ 5。玩这样的游戏你愿意支付多少钱？如果赢 \$ 500 而不是 \$ 5 又如何？
6. 推广购物篮分析的置信度和支持度公式，计算 k -依赖，即 $P(Y | X_1, \dots, X_k)$ 。
7. 证明：把一项由前件移到后件置信度不会增加： $\text{confidence}(ABC \rightarrow D) \geq \text{confidence}(AB \rightarrow CD)$ 。
8. 在购物篮分析中，如果对于每件售出的商品我们还有一个数，该数指出顾客喜爱该商品的程度，例如，在 0 到 10 这个范围内，你怎么能利用这一附加信息计算把哪种商品推荐给顾客？
9. 给出事务数据的例子，其中对于规则 $X \rightarrow Y$,
- (a) 支持度和置信度都高。
 - (b) 支持度高而置信度低。
 - (c) 支持度低而置信度高。
 - (d) 支持度和置信度都低。

3.9 参考文献

- Agrawal, R., H. Mannila, R. Srikant, H. Toivonen, and A. Verkamo. 1996. "Fast Discovery of Association Rules." In *Advances in Knowledge Discovery and Data Mining*, ed. U. M. Fayyad, G. Piatetsky-Shapiro, P. Smyth, and R. Uthurusamy, 307-328. Cambridge, MA: MIT Press.
- Duda, R. O., P. E. Hart, and D. G. Stork. 2001. *Pattern Classification*, 2nd ed. New York: Wiley.
- Li, J. 2006. "On Optimal Rule Discovery." *IEEE Transactions on Knowledge and Data Discovery* 18: 460-471.
- Newman, J. R., ed. 1988. *The World of Mathematics*. Redmond, WA: Tempus.
- Omicinski, E. R. 2003. "Alternative Interest Measures for Mining Associations in Databases." *IEEE Transactions on Knowledge and Data Discovery* 15: 57-69.
- Russell, S., and P. Norvig. 1995. *Artificial Intelligence: A Modern Approach*. New York: Prentice Hall.
- Shafer, G., and J. Pearl, eds. 1990. *Readings in Uncertain Reasoning*. San Mateo, CA: Morgan Kaufmann.
- Zhang, C., and S. Zhang. 2002. *Association Rule Mining: Models and Algorithms*. New York: Springer.

参数方法

前面, 我们讨论了在使用概率对不确定性建模时, 如何做出最优决策。现在, 考虑如何从给定的训练集估计这些概率。我们从分类和回归的参数方法开始。在后面的章节中, 将讨论半参数和非参数方法, 介绍用于权衡模型复杂度和经验误差的模型选择方法。

4.1 引言

统计量 (statistic) 是从给定样本中计算的任何值。在统计推断中, 我们使用样本提供的信息进行决策。第一种方法是参数方法, 这里假设样本取自服从已知模型的某个分布, 例如高斯分布。参数方法的优点是, 模型定义在少量参数 (例如均值、方差), 即分布的有效统计量上。一旦从样本中估计出这些参数, 就知道了整个分布。我们从给定的样本估计分布的参数, 把这些估计放到假设的模型中, 并得到估计的分布, 然后用它进行决策。我们用来估计分布参数的方法是最大似然估计。我们还将讨论贝叶斯估计, 随着计算能力的提高, 贝叶斯估计日趋流行。

我们从密度估计 (density estimation) 开始。密度估计是估计 $p(x)$ 的一般情况。我们使用密度估计进行分类, 其中估计的密度是能够计算后验概率 $P(C_i | x)$ 并做决策的类密度 $p(x | C_i)$ 和 $P(C_i)$ 。然后, 我们讨论回归, 其中估计的密度是 $p(y | x)$ 。本章, x 是一维的, 因此密度是一元的。在第5章中, 我们将推广到多元情况。

60
3
61

4.2 最大似然估计

假定我们有一个独立同分布 (iid) 样本 $\mathcal{X} = \{x^t\}_{t=1}^N$ 。我们假设 x^t 是从某个定义在参数 θ 上的已知概率密度族 $p(x | \theta)$ 中抽取的实例:

$$x^t \sim p(x | \theta)$$

我们希望找出这样的 θ , 使得 x^t 尽可能像是从 $p(x | \theta)$ 抽取的。因为 x^t 是独立的, 给定参数 θ , 样本 \mathcal{X} 的似然 (likelihood) 是个体点似然的乘积:

$$l(\mathcal{X} | \theta) \equiv p(\mathcal{X} | \theta) = \prod_{t=1}^N p(x^t | \theta) \quad (4.1)$$

在最大似然估计 (maximum likelihood estimation) 中, 我们感兴趣的是找到这样的 θ , 使得 \mathcal{X} 最像是抽取的。因此, 我们寻找最大化样本似然的 θ , 该似然记作 $l(\theta | \mathcal{X})$ 。我们可以最大化该似然的对数, 而不改变它取最大值的数值。 $\log(\cdot)$ 把乘积转换为求和, 并且当假定某种密度 (例如包含指数) 时进一步简化计算量。对数似然 (log likelihood) 定义为

$$\mathcal{L}(\theta | \mathcal{X}) \equiv \log l(\theta | \mathcal{X}) = \sum_{t=1}^N \log p(x^t | \theta) \quad (4.2)$$

现在, 让我们来看我们感兴趣的实际应用中出现的一些分布。如果我们有两类问题, 我们就

62

用伯努利分布。当存在 $K > 2$ 个类时, 分布拓广为多项分布。高斯(正态)密度是一种最常使用的对类条件输入建模的密度。对于这三种分布, 我们讨论它们的参数的最大似然估计(MLE)。

4.2.1 伯努利密度

在伯努利分布中, 有两个结果: 事件要么发生, 要么不发生。例如, 实例是类的正例, 或者不是。事件发生, 伯努利随机变量 X 以概率 p 取值 1, 事件不发生的概率为 $1 - p$, 并用 X 取值 0 表示。这表示为

$$P(x) = p^x(1-p)^{1-x}, \quad x \in \{0, 1\} \quad (4.3)$$

期望值和方差可以用下式计算:

$$E[X] = \sum_x xp(x) = 1 \cdot p + 0 \cdot (1-p) = p$$

$$\text{Var}(X) = \sum_x (x - E[X])^2 p(x) = p(1-p)$$

p 是唯一参数, 并且给定独立同分布样本 $\mathcal{X} = \{x^t\}_{t=1}^N$, 其中 $x^t \in \{0, 1\}$, 我们希望计算 p 的估计 \hat{p}_i 。对数似然是

$$\begin{aligned} \mathcal{L}(p|\mathcal{X}) &= \log \prod_{t=1}^N p^{(x^t)} (1-p)^{(1-x^t)} \\ &= \sum_t x^t \log p + \left(N - \sum_t x^t\right) \log(1-p) \end{aligned}$$

通过求解 $d\mathcal{L}/dp = 0$ 可以找出最大化该对数似然的 \hat{p}_i , p 上带帽表示它是 p 的一个估计。

$$\hat{p}_i = \frac{\sum_t x^t}{N} \quad (4.4)$$

p 的估计是事件发生的次数与试验次数的比值。注意, 如果 X 是参数为 p 的伯努利变量, 则 $E[X] = p$, 并且作为期望, 均值的最大似然估计是样本的平均值。

注意, 该估计是样本的函数, 并且也是一个随机变量; 给定从相同的 $p(x)$ 中抽取的不同的 x_i , 我们可以谈论 \hat{p}_i 分布。例如, \hat{p}_i 分布的方差预料随 N 的增加而减少; 随着样本增大, 它们(从而它们的平均值)变得更相似。

63

4.2.2 多项密度

以伯努利分布的推广为例, 其中随机事件的结果不是两种状态, 而是 K 种互斥、穷举状态之一(例如类), 其中每种状态出现的概率为 p_i , 其中 $\sum_{i=1}^K p_i = 1$ 。设 x_1, x_2, \dots, x_K 是指示变量, 其中当输出为状态 i 时 x_i 为 1, 否则为 0。

$$P(x_1, x_2, \dots, x_K) = \prod_{i=1}^K p_i^{x_i} \quad (4.5)$$

假定我们做 N 次这样的独立试验, 结果为 $\mathcal{X} = \{x^t\}_{t=1}^N$, 其中,

$$x_i^t = \begin{cases} 1 & \text{如果试验 } t \text{ 选择状态 } i \\ 0 & \text{否则} \end{cases}$$

其中 $\sum_i x_i^t = 1$ 。 p_i 的最大似然估计是

$$\hat{p}_i = \frac{\sum_t x_i^t}{N} \quad (4.6)$$

状态 i 的概率估计是结果为状态 i 的试验次数与试验总次数的比值。有两种方法可以获得这个估计：如果 x_i 是 0/1，则可以认为它们是 K 个独立的伯努利试验。或者，我们可以写出对数似然并找出最大化它的 p_i (满足条件 $\sum_i p_i = 1$)。

4.2.3 高斯(正态)密度

X 是均值为 μ 方差为 σ^2 的高斯(正态)分布，记作 $\mathcal{N}(\mu, \sigma^2)$ ，如果它的密度函数为

$$p(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{(x-\mu)^2}{2\sigma^2}\right], \quad -\infty < x < \infty \quad (4.7)$$

给定样本 $\mathcal{X} = \{x^t\}_{t=1}^N$ ，其中 $x^t \sim \mathcal{N}(\mu, \sigma^2)$ ，高斯样本的对数似然为

$$\mathcal{L}(\mu, \sigma | \mathcal{X}) = -\frac{N}{2} \log(2\pi) - N \log \sigma - \frac{\sum_t (x^t - \mu)^2}{2\sigma^2}$$

通过求该对数似然的偏导数并令它们等于零，求得最大似然估计为

$$\begin{aligned} m &= \frac{\sum_t x^t}{N} \\ s^2 &= \frac{\sum_t (x^t - m)^2}{N} \end{aligned} \quad (4.8)$$

64

我们根据通常的约定，用希腊字母表示总体参数，用罗马字母表示它们的样本估计。有时，帽(抑扬符号)也用来表示估计，例如 $\hat{\mu}$ 。

4.3 评价估计：偏倚和方差

令 \mathcal{X} 是取自参数 θ 指定的总体上的样本，并令 $d = d(\mathcal{X})$ 是 θ 的一个估计。为了评估该估计的质量，我们可以度量它与 θ 有多大不同，即 $(d(\mathcal{X}) - \theta)^2$ 。但是因为它是一个随机变量(它依赖于样本)，我们需要对它在可能的 \mathcal{X} 上取平均值，并考虑 $r(d, \theta)$ ，它是估计 d 的均方误差(mean square error)，定义为

$$r(d, \theta) = E[(d(\mathcal{X}) - \theta)^2] \quad (4.9)$$

估计的偏倚(bias)是

$$b_\theta(d) = E[d(\mathcal{X})] - \theta \quad (4.10)$$

如果对所有的 θ 值都有 $b_\theta(d) = 0$ ，则 d 是 θ 的无偏估计(unbiased estimator)。例如，如果 x^t 是从均值为 μ 的密度抽取出的，则样本平均值 m 是均值 μ 的一个无偏估计，因为

$$E(m) = E\left[\frac{\sum_t x^t}{N}\right] = \frac{1}{N} \sum_t E[x^t] = \frac{N\mu}{N} = \mu$$

这就意味着虽然在一个特定样本上， m 却可能与 μ 不同，但是如果我们取许多这样的样本 \mathcal{X}_i ，并且估计许多 $m_i = m(\mathcal{X}_i)$ ，随着样本的增加，它们的平均值将逼近 μ 。 m 也是一个一致估计，也就是说，当 $N \rightarrow \infty$ 时， $\text{Var}(m) \rightarrow 0$ 。

$$\text{Var}(m) = \text{Var}\left(\frac{\sum_t x^t}{N}\right) = \frac{1}{N^2} \sum_t \text{Var}(x^t) = \frac{N\sigma^2}{N^2} = \frac{\sigma^2}{N}$$

随着样本中的点数 N 的增大， m 对 μ 的偏离变小。现在，让我们来检查 σ^2 的最大似然估计 s^2 ：

$$s^2 = \frac{\sum_i (x^i - m)^2}{N} = \frac{\sum_i (x^i)^2 - Nm^2}{N}$$

$$E[s^2] = \frac{\sum_i E[(x^i)^2] - N \cdot E[m^2]}{N}$$

给定 $\text{Var}(X) = E[X^2] - E[X]^2$, 我们得到 $E[X^2] = \text{Var}(X) + E[X]^2$, 并且

65

$$E[(x^i)^2] = \sigma^2 + \mu^2, \quad E[m^2] = \sigma^2/N + \mu^2$$

于是, 我们有

$$E[s^2] = \frac{N(\sigma^2 + \mu^2) - N(\sigma^2/N + \mu^2)}{N} = \left(\frac{N-1}{N}\right)\sigma^2 \neq \sigma^2$$

上式说明 s^2 是 σ^2 的有偏估计。 $(N/(N-1))s^2$ 是一个无偏估计。然而, 当 N 很大时, 差别可以忽略。这是一个渐近无偏估计 (asymptotically unbiased estimator) 的例子, 它的偏倚随着 N 趋向无穷而趋向于 0。

均方误差可以重新改写如下 (d 是 $d(X)$ 的缩写):

$$\begin{aligned} r(d, \theta) &= E[(d - \theta)^2] \\ &= E[(d - E[d] + E[d] - \theta)^2] \\ &= E[(d - E[d])^2 + (E[d] - \theta)^2 + 2(E[d] - \theta)(d - E[d])] \\ &= E[(d - E[d])^2] + E[(E[d] - \theta)^2] + 2E[(E[d] - \theta)(d - E[d])] \\ &= E[(d - E[d])^2] + (E[d] - \theta)^2 + 2(E[d] - \theta)E[d - E[d]] \\ &= \underbrace{E[(d - E[d])^2]}_{\text{方差}} + \underbrace{(E[d] - \theta)^2}_{\text{偏倚}^2} \end{aligned} \quad (4.11)$$

最后两式相等是因为 $E[d]$ 是常数, 因此 $E[d] - \theta$ 也是一个常数, 并且因为 $E[d - E[d]] = E[d] - E[d] = 0$ 。在式 (4.11) 中, 第一项是方差 (variance), 度量在平均情况下 d_i 在期望值附近的变化程度 (从一个数据集到另一个数据集); 而第二项是偏倚 (bias), 度量期望值偏离正确值 θ 的程度 (参见图 4-1)。于是, 我们把误差写成方差和偏倚的平方之和:

$$r(d, \theta) = \text{Var}(d) + (b_\theta(d))^2 \quad (4.12)$$

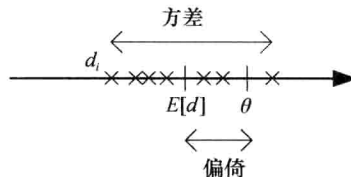


图 4-1 θ 是需要估计的参数。 d_i 是在不同样本上的多个估计 (用“×”表示)。偏倚是 d 的期望值与 θ 之差。方差是 d_i 在期望值周围的散布程度。我们希望它们两个都很小

4.4 贝叶斯估计

有时, 在看到样本之前, 我们 (或应用领域专家) 可能会有一些关于参数 θ 可能取值的范围的先验 (prior) 信息。这些信息是非常有用的, 也应当利用起来, 尤其是样本较小时。

这些先验信息不会告诉我们参数的确切值(否则我们就不需要该样本),并且我们通过把 θ 看作是一个随机变量并为其定义先验密度 $p(\theta)$ 来对这种不确定性建模。例如,假设我们知道 θ 接近正态分布,并且 θ 介于 5 到 9 之间,7 左右对称的置信度为 90%。于是,可以把 $p(\theta)$ 写成均值为 7 的正态分布,并且因为

$$P\left\{-1.64 < \frac{\theta - \mu}{\sigma} < 1.64\right\} = 0.9$$

$$P\{\mu - 1.64\sigma < \theta < \mu + 1.64\sigma\} = 0.9$$

我们取 $1.64\sigma = 2$, 并且使用 $\sigma = 2/1.64$ 。这样,我们就可以假定 $p(\theta) \sim \mathcal{N}(7, (2/1.64)^2)$ 。

先验密度(prior density) $p(\theta)$ 告诉我们在看到样本之前 θ 的可能取值。我们把它和样本数据告诉我们的(即似然密度 $p(X|\theta)$)结合起来,利用贝叶斯规则,得到 θ 的后验密度(posterior density),它告诉我们看到样本之后 θ 的可能取值:

$$p(\theta|X) = \frac{p(X|\theta)p(\theta)}{p(X)} = \frac{p(X|\theta)p(\theta)}{\int p(X|\theta')p(\theta')d\theta'} \quad (4.13)$$

为了估计 x 上的密度,我们有

$$\begin{aligned} p(x|X) &= \int p(x, \theta|X) d\theta \\ &= \int p(x|\theta, X) p(\theta|X) d\theta \\ &= \int p(x|\theta) p(\theta|X) d\theta \end{aligned}$$

$p(x|\theta, X) = p(x|\theta)$, 因为只要我们知道有效统计量 θ , 我们就知道关于分布的一切。这样,我们在所有 θ 的值的预测上取平均值,用它们的概率加权。如果我们像在回归中一样,以 $y = g(x|\theta)$ 的形式做预测,则有

$$y = \int g(x|\theta) p(\theta|X) d\theta$$

除非后验具有很好的形式,否则求这个积分可能非常困难。当求整个积分不可行时,我们把它缩减到单个点。如果我们可以假定 $p(\theta|X)$ 在它的众数周围有一个窄的峰值,则使用最大后验(maximum a posteriori, MAP)估计将使得计算比较容易:

$$\theta_{MAP} = \arg \max_{\theta} p(\theta|X) \quad (4.14)$$

这样,用单个点取代整个密度,回避积分并且使用

$$\begin{aligned} p(x|X) &= p(x|\theta_{MAP}) \\ y_{MAP} &= g(x|\theta_{MAP}) \end{aligned}$$

如果我们没有更重要的理由偏爱 θ 的某些值,则先验密度是扁平的,后验将与似然 $p(X|\theta)$ 有同样的形式,并且 MAP 估计将等价于最大似然估计(参见 4.2 节),其中我们有

$$\theta_{ML} = \arg \max_{\theta} p(X|\theta) \quad (4.15)$$

另外一个可能的方法是贝叶斯估计(Bayes' estimator),它被定义为后验密度的期望值

$$\theta_{\text{Bayes}} = E[\theta | \mathcal{X}] = \int \theta p(\theta | \mathcal{X}) d\theta \quad (4.16)$$

取期望值的原因是随机变量的最佳估计是它的均值。假设 θ 是变量，我们想要用 $E[\theta] = \mu$ 预测。可以证明如果常数 c 是 θ 的估计，则

$$\begin{aligned} E[(\theta - c)^2] &= E[(\theta - \mu + \mu - c)^2] \\ &= E[(\theta - \mu)^2] + (\mu - c)^2 \end{aligned} \quad (4.17)$$

68

如果 c 取 μ ，它的值最小。在正态密度情况下，众数是期望值。这样，如果 $p(\theta | \mathcal{X})$ 是正态的，则 $\theta_{\text{Bayes}} = \theta_{\text{MAP}}$ 。

作为一个例子，我们假设 $x' \sim \mathcal{N}(\theta, \sigma_0^2)$ 并且 $\theta \sim \mathcal{N}(\mu, \sigma^2)$ ，其中 μ ， σ 和 σ_0^2 已知：

$$\begin{aligned} p(\mathcal{X} | \theta) &= \frac{1}{(2\pi)^{N/2} \sigma_0^N} \exp \left[-\frac{\sum_i (x'_i - \theta)^2}{2\sigma_0^2} \right] \\ p(\theta) &= \frac{1}{\sqrt{2\pi}\sigma} \exp \left[-\frac{(\theta - \mu)^2}{2\sigma^2} \right] \end{aligned}$$

可以证明 $p(\theta | \mathcal{X})$ 是正态的，满足

$$E[\theta | \mathcal{X}] = \frac{N/\sigma_0^2}{N/\sigma_0^2 + 1/\sigma^2} m + \frac{1/\sigma^2}{N/\sigma_0^2 + 1/\sigma^2} \mu \quad (4.18)$$

因此，贝叶斯估计是先验均值 μ 和样本均值 m 的加权平均值，权重与它们的方差成反比。利用样本提供的更多的信息，随着样本规模 N 的增加，贝叶斯估计逼近样本的平均值。当 σ^2 较小时，即当我们关于 θ 正确值具有较少的先验不确定性时，或者当 N 较小时，我们的先验猜测 μ 具有较好的效果。

注意，MAP 和贝叶斯估计都把整个后验密度归约到单个点并且丢失信息，除非后验是单模的并且在这些点周围有一个窄峰。随着计算费用降低，一个可能性是使用蒙特卡洛方法，它可以从后验密度产生样本 (Andrieu 等 2003)。还有一些近似方法可以用来计算整个积分。我们将在第 14 章更详细地讨论贝叶斯估计。

4.5 参数分类

我们在第 3 章看到，使用贝叶斯规则，我们可以把类 C_i 的后验概率写为

$$P(C_i | x) = \frac{p(x | C_i) P(C_i)}{p(x)} = \frac{p(x | C_i) P(C_i)}{\sum_{k=1}^K p(x | C_k) P(C_k)} \quad (4.19)$$

并使用判别式函数

$$g_i(x) = p(x | C_i) P(C_i)$$

或者等价地

$$g_i(x) = \log p(x | C_i) + \log P(C_i) \quad (4.20)$$

如果我们可以假设 $p(x | C_i)$ 是高斯的，

$$p(x | C_i) = \frac{1}{\sqrt{2\pi}\sigma_i} \exp \left[-\frac{(x - \mu_i)^2}{2\sigma_i^2} \right] \quad (4.21)$$

则式(4.20)变成

$$g_i(x) = -\frac{1}{2} \log 2\pi - \log \sigma_i - \frac{(x - \mu_i)^2}{2\sigma_i^2} + \log P(C_i) \quad (4.22)$$

让我们看一个例子：假设一个汽车销售 K 种不同的汽车，为了简单起见，我们假定唯一影响顾客购买的因素是他们的年收入，用 x 表示。于是， $P(C_i)$ 是购买类型 i 汽车的顾客所占的比例。如果顾客的年收入分布可以用一个高斯分布近似，则购买类型 i 汽车的顾客其收入为 x 的概率 $p(x | C_i)$ 服从分布 $\mathcal{N}(\mu_i, \sigma_i^2)$ ，其中 μ_i 是这类顾客年收入的均值， σ_i^2 是他们的年收入的方差。

当我们不知道 $P(C_i)$ 和 $p(x | C_i)$ 时，我们从样本估计它们并把它们的估计插入判别式，得到判别函数的估计。我们有样本

$$\mathcal{X} = \{x^t, r^t\}_{i=1}^N \quad (4.23)$$

其中 $x \in \mathcal{R}$ 是一维的， $r \in \{0, 1\}^K$ 使得

$$r_i^t = \begin{cases} 1 & \text{如果 } x^t \in C_i \\ 0 & \text{如果 } x^t \in C_k, k \neq i \end{cases} \quad (4.24)$$

对于每一个类，均值和方差的估计是(依赖于 4.8 式)

$$m_i = \frac{\sum_t x^t r_i^t}{\sum_t r_i^t} \quad (4.25)$$

$$s_i^2 = \frac{\sum_t (x^t - m_i)^2 r_i^t}{\sum_t r_i^t} \quad (4.26)$$

而先验的估计是(依赖于 4.6 式)

$$\hat{p}_i(C_i) = \frac{\sum_t r_i^t}{N} \quad (4.27) \quad \boxed{70}$$

把这些估计代入式(4.22)，得到

$$g_i(x) = -\frac{1}{2} \log 2\pi - \log s_i - \frac{(x - m_i)^2}{2s_i^2} + \log \hat{P}(C_i) \quad (4.28)$$

第一项是常数，可以去掉，因为它在所有的 $g_i(x)$ 中都是一样的。如果这些先验相等，则最后一项也可以去掉。如果我们进一步假设方差都相等，则上式可以写为

$$g_i(x) = -(x - m_i)^2 \quad (4.29)$$

因此我们把 x 指派到具有最近均值的类：

$$\text{选择 } C_i, \text{ 如果 } |x - m_i| = \min_k |x - m_k|$$

对于两个相邻的类，两个均值之间的中点是决策阈值(参见图 4-2)。

$$g_1(x) = g_2(x) \quad \boxed{71}$$

$$(x - m_1)^2 = (x - m_2)^2$$

$$x = \frac{m_1 + m_2}{2}$$

当方差不相同时有两个阈值(参见图 4-3)，它们都容易计算(参见习题 4)。如果先验概率不同，则具有向不太可能的类的均值移动决策阈值的效果。

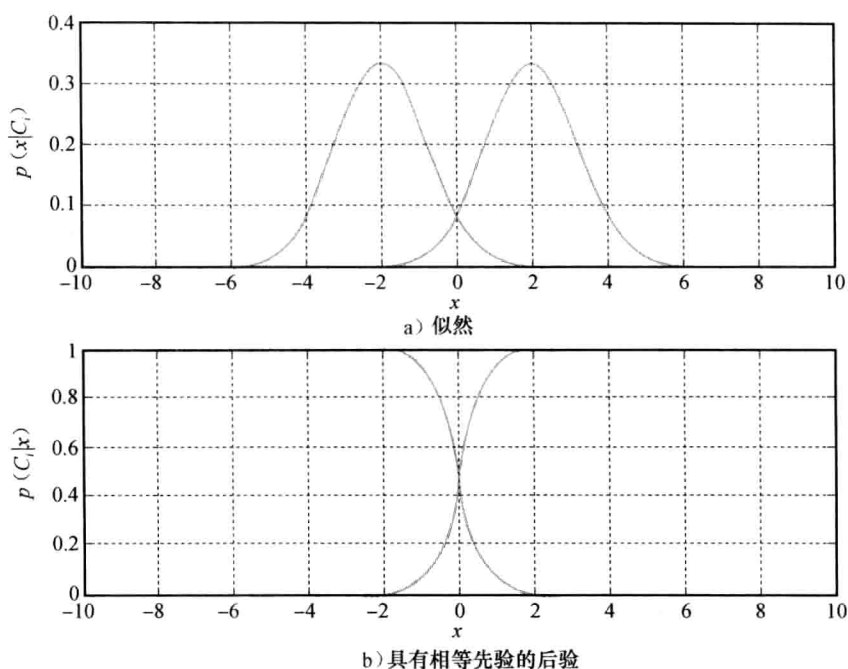


图4-2 输入是一维的，似然函数 a) 和具有相等先验的后验的两个类 b)。方差相等并且后验相交于一点，该点是决策阈值

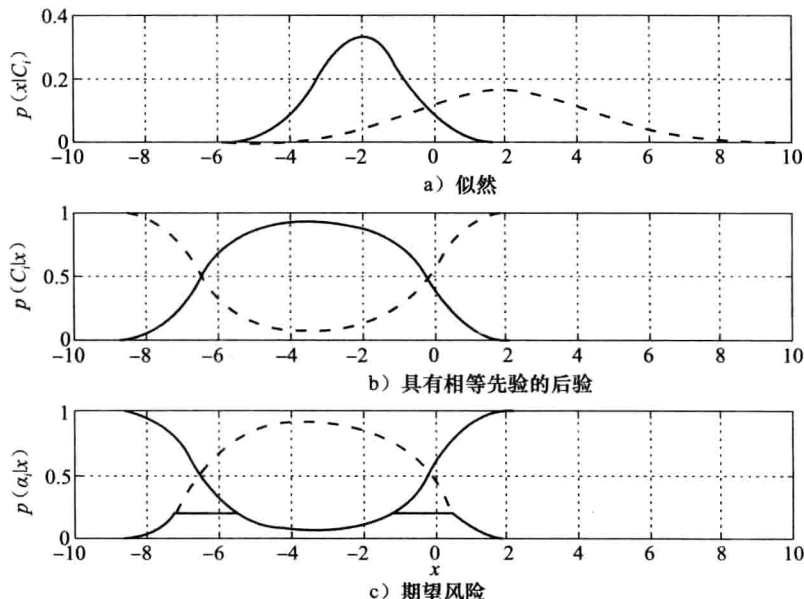


图4-3 输入是一维的，似然函数 a) 和具有相等先验的后验的两个类 b)。方差不相等并且后验在两个点上相交。c) 中对两个类和 $\lambda = 0.2$ 的拒绝 (参见 3.3 节) 显示了期望风险

这里，我们对参数使用最大似然估计。但是，如果有一些关于它们的先验信息 (如均值)，则我们可以利用 μ_i 上的先验，使用 $p(x | C_i)$ 的贝叶斯估计。

这里必须注意：当 x 连续时，我们不要急于对 $p(x | C_i)$ 使用高斯密度。如果密度函数不

是高斯的，则分类算法(即阈值点)将会出错。在统计学文献中，存在检查正态性的检验，并且这样的检验应该在假定正态分布之前使用。在一维数据的情况下，最简单的检验是绘制直方图并观察密度是否是钟形的，即是否是单峰并且围绕中心对称。

这是基于似然(likelihood-based approach)的分类方法，其中我们使用数据估计密度，使用贝叶斯规则计算后验密度，然后得到判别式。在稍后的章节中，我们讨论基于判别式的方法，那里我们将绕过密度估计而直接估计判别式。

4.6 回归

在回归中，我们喜欢将数值输出写成输入的函数。数值输出称为因变量(dependent variable)，函数的输入称为自变量(independent variable)。我们假定数值输出是输入的确定性函数与随机噪声的和：

$$r = f(x) + \varepsilon$$

其中 $f(x)$ 是未知函数，我们将用定义在参数 θ 的集合上的估计 $g(x|\theta)$ 来近似它。如果我们假设 ε 服从均值为0，方差为 σ^2 的高斯分布，即 $\varepsilon \sim \mathcal{N}(0, \sigma^2)$ ，并且用我们的估计 $g(\cdot)$ 取代未知函数 $f(\cdot)$ ，则我们有(参见图4-4)

$$p(r|x) \sim \mathcal{N}(g(x|\theta), \sigma^2) \quad (4.30)$$

我们再一次使用最大似然来学习参数 θ 。训练集中的 (x^i, r^i) 对偶取自未知的联合概率密度 $p(x, r)$ ，可以写作

$$p(x, r) = p(r|x)p(x)$$

$p(r|x)$ 是在给定输入下输出的概率，而 $p(x)$ 是输入的密度。给定 iid 样本 $\mathcal{X} = \{x^i, r^i\}_{i=1}^N$ ，对数似然是

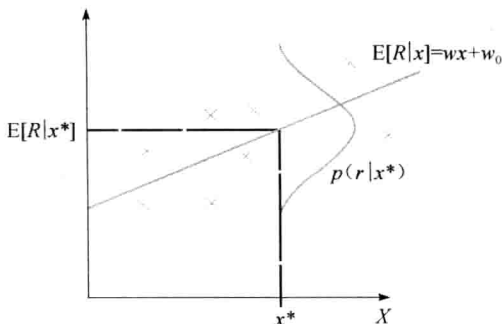


图4-4 回归假定0均值的高斯噪声加到模型中，这里模型是线性的

$$\mathcal{L}(\theta|\mathcal{X}) = \log \prod_{i=1}^N p(x^i, r^i) = \log \prod_{i=1}^N p(r^i|x^i) + \log \prod_{i=1}^N p(x^i)$$

我们可以忽略第二项，因为它不依赖于我们的估计。于是，我们有

$$\begin{aligned} \mathcal{L}(\theta|\mathcal{X}) &= \log \prod_{i=1}^N \frac{1}{\sqrt{2\pi}\sigma} \exp \left[-\frac{[r^i - g(x^i|\theta)]^2}{2\sigma^2} \right] \\ &= \log \left(\frac{1}{\sqrt{2\pi}\sigma} \right)^N \exp \left[-\frac{1}{2\sigma^2} \sum_{i=1}^N [r^i - g(x^i|\theta)]^2 \right] \\ &= -N \log(\sqrt{2\pi}\sigma) - \frac{1}{2\sigma^2} \sum_{i=1}^N [r^i - g(x^i|\theta)]^2 \end{aligned} \quad (4.31)$$

第一项独立于参数 θ ，可以去掉，因子 $1/\sigma^2$ 也可以去掉。最大化上式等价于最小化

$$E(\theta|\mathcal{X}) = \frac{1}{2} \sum_{i=1}^N [r^i - g(x^i|\theta)]^2 \quad (4.32)$$

它是我们最经常使用误差函数，而最小化它的 θ 叫做最小二乘估计(least squares estimate)。这是统计学经常做的一个变换：当似然 l 包含指数时，取代最小化 l ，我们定义一个误差函

数 $E = -\log l$, 并最小化它。

在线性回归(linear regression)中, 我们有线性模型

$$g(x^t | w_1, w_0) = w_1 x^t + w_0$$

取误差平方和的导数(4.32式), 我们得到两个未知数的两个方程

$$\begin{aligned}\sum_t r^t &= Nw_0 + w_1 \sum_t x^t \\ \sum_t r^t x^t &= w_0 \sum_t x^t + w_1 \sum_t (x^t)^2\end{aligned}$$

它们可以被写成向量矩阵的形式 $\mathbf{A}\mathbf{w} = \mathbf{y}$, 其中

$$\mathbf{A} = \begin{bmatrix} N & \sum_t x^t \\ \sum_t x^t & \sum_t (x^t)^2 \end{bmatrix}, \quad \mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} \sum_t r^t \\ \sum_t r^t x^t \end{bmatrix}$$

并且可以解得 $\mathbf{w} = \mathbf{A}^{-1}\mathbf{y}$ 。

在多项式回归(polynomial regression)的一般情况下, 该模型是 k 次多项式

$$g(x^t | w_k, \dots, w_2, w_1, w_0) = w_k (x^t)^k + \dots + w_2 (x^t)^2 + w_1 x^t + w_0$$

这个模型关于它的参数是线性的, 取它的导数, 我们可以得到 $k+1$ 个未知数的 $k+1$ 个方程, 可以写做向量矩阵的形式 $\mathbf{A}\mathbf{w} = \mathbf{y}$, 其中我们有

$$\mathbf{A} = \begin{bmatrix} N & \sum_t x^t & \sum_t (x^t)^2 & \dots & \sum_t (x^t)^k \\ \sum_t x^t & \sum_t (x^t)^2 & \sum_t (x^t)^3 & \dots & \sum_t (x^t)^{k+1} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \sum_t (x^t)^k & \sum_t (x^t)^{k+1} & \sum_t (x^t)^{k+2} & \dots & \sum_t (x^t)^{2k} \end{bmatrix}$$

$$\mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ \vdots \\ w_k \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} \sum_t r^t \\ \sum_t r^t x^t \\ \sum_t r^t (x^t)^2 \\ \vdots \\ \sum_t r^t (x^t)^k \end{bmatrix}$$

记作 $\mathbf{A} = \mathbf{D}^T \mathbf{D}$ 和 $\mathbf{y} = \mathbf{D}^T \mathbf{r}$, 其中

$$\mathbf{D} = \begin{bmatrix} 1 & x^1 & (x^1)^2 & \dots & (x^1)^k \\ 1 & x^2 & (x^2)^2 & \dots & (x^2)^k \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x^N & (x^N)^2 & \dots & (x^N)^k \end{bmatrix}, \quad \mathbf{r} = \begin{bmatrix} r^1 \\ r^2 \\ \vdots \\ r^N \end{bmatrix}$$

然后, 我们可以求解参数, 得到

$$\mathbf{w} = (\mathbf{D}^T \mathbf{D})^{-1} \mathbf{D}^T \mathbf{r} \quad (4.33)$$

假设高斯分布误差, 则最大化似然对应于最小化误差的平方和。另外一个度量是相对平方误差(relative square error, RSE)

$$E_{\text{RSE}} = \frac{\sum_t [r^t - g(x^t | \theta)]^2}{\sum_t (r^t - \bar{r})^2} \quad (4.34)$$

如果 E_{RSE} 接近于 1, 则我们的预测与用平均值的预测一样好; 当它接近于 0, 我们得到

更好的拟合。如果 E_{RSE} 接近于 1，则意味着使用基于输入 x 的模型不比没有 x 的情况下使用平均值作为估计器更好；如果 E_{RSE} 接近于 0，则输入 x 是有用的。

为了检查回归是否实现很好的拟合，一个度量是决定系数(coefficient of determination)

$$R^2 = 1 - E_{RSE}$$

并且为了对回归是有用的，我们要求 R^2 接近于 1。

记住，为了最好地泛化，我们应该调整学习模型的复杂度，以适应数据的复杂度。在多项式回归中，复杂度参数是拟合多项式的阶，因此需要找到一种选择最佳阶数的方法，能够最小化泛化误差。也就是说，找到一种方法，调整模型的复杂度使其最佳拟合数据所固有的函数复杂度。

75

4.7 调整模型的复杂度：偏倚/方差两难选择

我们假设样本 $\mathcal{X} = \{x^t, r^t\}$ 取自未知的联合概率密度 $p(x, r)$ 。使用这个样本，我们构建估计 $g(\cdot)$ 。 x 上(联合密度上)的期望平方误差可以表示为(用 4.17 式)

$$E[(r - g(x))^2 | x] = \underbrace{E[(r - E[r|x])^2 | x]}_{\text{噪声}} + \underbrace{(E[r|x] - g(x))^2}_{\text{误差的平方}} \quad (4.35)$$

右边的第一项是给定 x 时 r 的方差；它不依赖于 $g(\cdot)$ 或 \mathcal{X} 。它是添加噪声的方差 σ^2 。它是误差的一部分，无论我们使用什么估计方法，都不可能消除它。第二项量化 $g(x)$ 偏离回归函数 $E[r|x]$ 的程度。它确实依赖于估计方法和训练集。对一个样本来说， $g(x)$ 也许是一个非常好的拟合；而对其他样本，它可能是很差的拟合。为了评价一个估计 $g(\cdot)$ 的好坏程度，我们对可能的数据集进行平均。

76

期望值(样本 \mathcal{X} 上的平均，所有的样本大小均为 N 并从相同联合密度 $p(x, r)$ 抽取)是(使用 4.11 式)

$$E_x[(E[r|x] - g(x))^2 | x] = \underbrace{(E[r|x] - E_x[g(x)])^2}_{\text{偏倚}} + \underbrace{E_x[(g(x) - E_x[g(x)])^2]}_{\text{方差}} \quad (4.36)$$

正如我们前面所讨论的，偏倚度量不考虑样本变化的影响时 $g(x)$ 的错误程度；方差度量当样本变化时 $g(x)$ 在期望值 $E[g(x)]$ 附近波动的程度。

让我们看一个例子：为了估计偏倚和方差，我们由某个带噪声的已知的 $f(\cdot)$ 产生一组数据集 $\mathcal{X}_i = \{x_i^t, r_i^t\}$ ， $i = 1, \dots, M$ ，利用每个数据集形成一个估计 $g_i(\cdot)$ ，并计算偏倚和方差。注意，在现实生活中，我们不能够这么做，因为我们不知道 $f(\cdot)$ ，也不知道所添加噪声的参数。于是， $E[g(x)]$ 用 $g_i(\cdot)$ 上的平均来估计：

$$\bar{g}(x) = \frac{1}{M} \sum_{i=1}^M g_i(x)$$

偏倚和方差的估计是

$$\text{Bias}^2(g) = \frac{1}{N} \sum_t [\bar{g}(x^t) - f(x^t)]^2$$

$$\text{Variance}(g) = \frac{1}{NM} \sum_t \sum_i [g_i(x^t) - \bar{g}(x^t)]^2$$

让我们看几个不同复杂度的模型：最简单的是常数拟合

$$g_i(x) = 2$$

它没有方差，因为我们没有使用数据，并且所有的 $g_i(x)$ 都是相同的。但是偏倚很高，除非对于所有的 x ， $f(x)$ 值都接近于 2。如果我们取样本中 r^i 的平均值

$$g_i(x) = \sum_t r_t^i / N$$

而不是常数 2，则就减少偏倚，因为我们预料在通常情况下，平均值是比常数更好的估计。
[77] 但是，这增加了方差，因为不同的样本 x_i 将有不同的平均值。通常，在这种情况下，偏倚的减少比方差的增加更大，而误差将会降低。

图 4-5 给出了多项式回归的情况下的一个例子。随着多项式阶的增加，数据集的较小的变化将导致拟合多项式的较大变化；因此方差增加。但是，复杂的模型可以更好地拟合潜在的函数；因此偏倚减少（参见图 4-6）。这称为偏倚/方差两难选择 (bias/variance dilemma)，并且不仅对于多项式回归，而且对于任何机器学习系统都存在这一问题 (Geman、Bienenstock 和 Doursat 1992)。为了减少偏倚，冒着具有高方差的危险，模型应当是柔性的。如果方差保持较低，则我们可能不能很好地拟合数据，并且具有较高的偏倚。最佳模型是最好地权衡偏倚和方差的模型。

如果有偏倚，这指明我们的模型类不包含解，这是欠拟合 (underfitting)。如果有方差，则模型类过于一般，并且也学习噪声，这是过拟合 (overfitting)。如果 $g(\cdot)$ 是与 $f(\cdot)$ 同样的假设类 (例如同阶多项式)，则我们有一个无偏估计，并且估计的偏倚随着模型数量的增加而减小。这表明选择正确模型的误差降低效果 (在第 2 章，我们称之为归纳偏倚——这两处“偏倚”的使用是不同的，但并非不相关)。对于方差，它同样依赖于训练集的大小；由于样本导致的可变性随着样本规模的增加而减少。总结一下，为了取得小的误差值，我们应该有合适的归纳偏倚 (在统计意义上取得小偏倚)，并且有一个足够大的数据集，使得模型的可变性能够受到数据的约束。

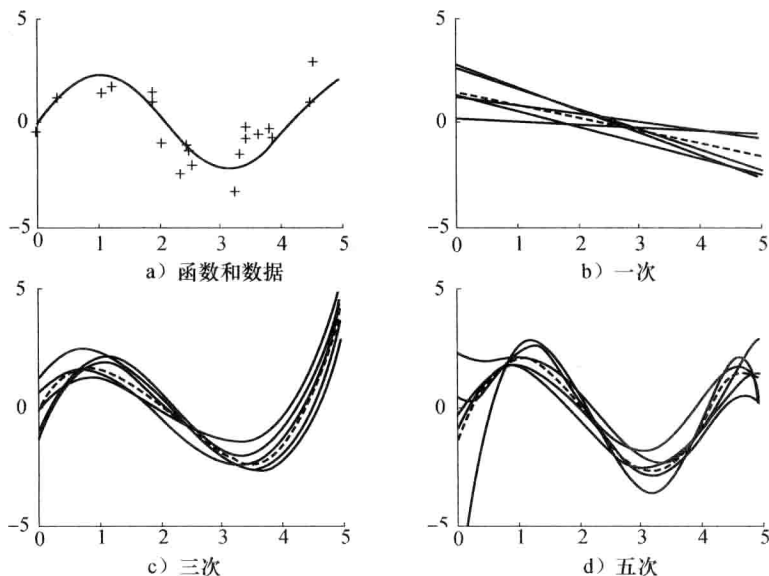


图 4-5 a) 函数 $f(x) = 2\sin(1.5x)$ 和一个从该函数抽样含噪声 ($\mathcal{N}(0, 1)$) 的数据集。抽取五个样本，每个包含 20 个实例。b)、c)、d) 分别是 5 个一次、三次和五次多项式拟合，即 $g_i(\cdot)$ 。对每种情况，虚线是 5 次拟合的平均 $\bar{g}(\cdot)$

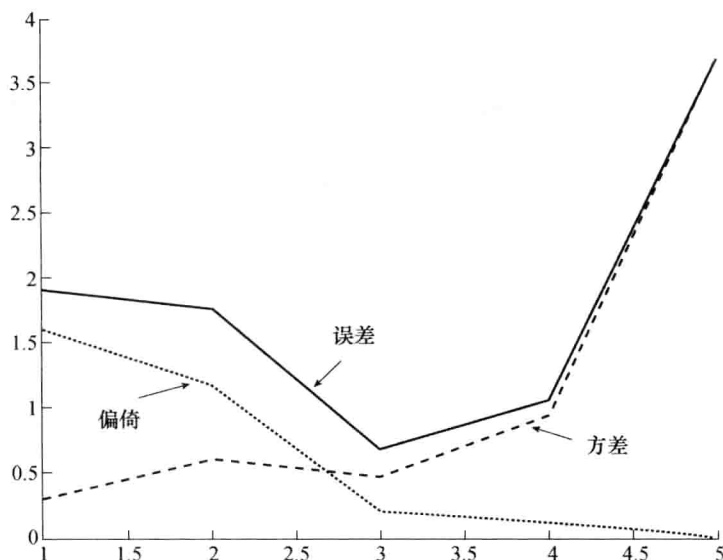


图 4-6 与图 4-5 同样的设置, 使用 100 个模型而不是 5 个, 从 1 到 5 次多项式的偏倚, 方差和误差。一次多项式具有最小方差, 5 次多项式具有最小偏倚。随着阶增加, 偏倚减小但是方差增加。3 次多项式具有最小误差

注意, 方差大时偏倚小, 表明 $\bar{g}(x)$ 是一个好的估计器。因此, 为了取得小误差, 我们可以采用大量的高方差模型, 并且用它们的平均值作为我们的估计。我们将在第 17 章讨论这种模型组合方法。

4.8 模型选择过程

有许多过程可以用来调整模型的复杂度。

在实践中, 我们用来发现最佳复杂度的方法是交叉验证 (cross-validation)。我们不能计算一个模型的偏倚和方差, 但是我们能够计算总误差。给定一个数据集, 我们把它分成两部分, 分别作为训练集和验证集, 在训练集上训练不同复杂度的候选模型, 在验证集上测试它们的误差。随着模型复杂度增加, 训练误差持续降低。在验证集的误差降低达到一定的复杂度水平之后, 停止降低或不再明显的降低, 如果有明显噪声的话甚至还会增加。这个“拐点”对应于最佳复杂度水平 (参见图 4-7)。

在现实生活中, 我们不能像图 4-6 中那样计算偏倚, 因而不能计算误差; 除了还包含噪声之外, 图 4-7 中的验证误差是一种估计: 即便我们有无偏的正确模型, 并且有足够大的数据集可以忽略方差, 也仍然会有非零验证误差。注意, 图 4-7 中的验证误差与图 4-6 中的误差不同, 不是“V 形”的, 因为前者使用更多训练数据, 并且我们知道使用更多数据可以限制方差。确实, 在图 4-5d) 中我们看到, 在有数据的地方, 即使 5 次多项式的表现也像 3 次多项式一样; 在只有少量数据的地方 (例如在两端), 5 次多项式不那么准确。

另外一个常用的方法是正则化 (regularization) (Breiman 1998a)。在这种方法中, 我们用一个增广误差函数, 记作

$$E' = \text{数据上的误差} + \lambda \cdot \text{模型复杂度} \quad (4.37)$$

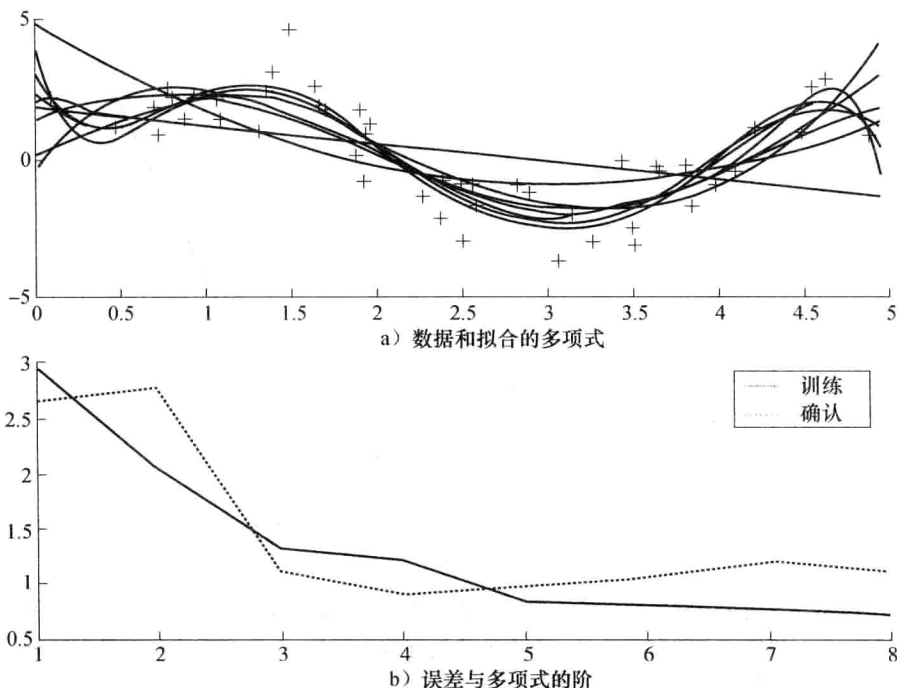


图 4-7 与图 4-5 同样的设置, 产生训练集和确认集(每个包含 50 个实例)。a) 训练数据和 1 到 8 阶拟合多项式。b) 训练和确认误差作为多项式阶的函数。“拐点”在 3

它的第二项用一个大的方差惩罚复杂模型, 其中 λ 给出了罚的权重。当我们最小化推广误差函数而不是数据上的误差时, 我们惩罚了复杂模型, 因此降低了方差。如果 λ 取得太大, 则只允许很简单的模型, 我们会冒着引进偏倚的危险。 λ 用交叉验证优化。

另一种看待式(4.37)的方法是把 E' 看作新的检验数据上的误差。右边的第一项是训练误差, 而第二项是估计训练与检验误差之间偏差的乐观 (optimism) 项 (Hastie, Tibshirani 和 Friedman 2001)。可以用诸如 Akaike 信息准则 (Akaike's information criterion, AIC) 和贝叶斯信息准则 (Bayesian information criterion, BIC) 等方法估计这个乐观项, 并把它加到训练误差上来估计检验误差, 而不必验证。这个乐观项的大小随输入数 d (这里是 $k+1$) 线性增加, 并且随训练集的大小 N 的增加而减少; 它还随添加噪声的方差 σ^2 (可以由低偏倚模型的误差估计) 增加。对于非线性模型, d 应该代之以参数的“有效”个数。

结构风险最小化 (structural risk minimization, SRM) (Vapnik 1995) 使用一个模型集, 按它们的复杂度排序。一个例子是阶递增的多项式。复杂度一般由自由参数的数量度量。VC 是另一种模型复杂度的度量。在式(4.37)中, 我们可以使用递减的 λ_i 来得到复杂度递增的模型集合。SRM 模型选择对应于寻找最简单并且在数据上的经验误差最小的模型。

最小描述长度 (minimum description length, MDL) (Rissanen 1978) 使用一种信息论度量。数据集的 Kolmogorov 复杂度定义为数据最短描述。如果数据简单, 它就有短的复杂度; 例如, 如果它是 0 的序列, 则我们可以只写 0 和序列的长度。如果数据完全随机的, 则我们不能有比数据自身更短的数据描述。如果一个模型对数据是合适的, 则它有一个好的数据拟合, 我们可以发送/存储数据描述而不是数据本身。在描述数据的所有模型中, 我们想要有

一个最简单的模型，那样它就可以有最短描述。这样，我们又一次要在模型的简单性和它解释数据能力之间权衡。

当我们有一些关于近似函数的合适类的先验知识时，我们使用贝叶斯模型选择 (Bayesian model selection)。这种先验知识被定义为模型的先验分布 $p(\text{模型})$ 。给定数据并假定一个模型，我们可以用贝叶斯规则计算 $p(\text{模型} | \text{数据})$ ：

$$p(\text{模型} | \text{数据}) = \frac{p(\text{数据} | \text{模型})p(\text{模型})}{p(\text{数据})} \quad (4.38)$$

$p(\text{模型} | \text{数据})$ 是给定我们关于模型的主观先验知识 (即 $p(\text{模型})$) 和数据提供的客观支持 (即 $p(\text{数据} | \text{模型})$)，模型的后验概率。我们可以选择具有最高后验概率的模型，或者用模型的后验概率加权，在所有模型上取平均。

如果我们取式 (4.38) 的对数，则得到

$$\log p(\text{模型} | \text{数据}) = \log p(\text{数据} | \text{模型}) + \log p(\text{模型}) - c \quad (4.39)$$

这具有式 (4.37) 的形式；数据的对数似然是训练误差，而先验的对数似然是罚项。例如，如果我们有一个回归模型，并使用先验 $p(\mathbf{w}) \sim \mathcal{N}(0, 1/\lambda)$ ，则 MAP 对应于最小的

$$E = \sum_i [r^i - g(x^i | \mathbf{w})]^2 + \lambda \sum_i w_i^2 \quad (4.40)$$

也就是说，我们寻找降低误差并且尽可能接近 0 的 w_i ，而我们希望它们接近 0 的理由是这样的话拟合的多项式会更平滑。随着多项式阶的增加，为了更好地拟合数据，函数将上下摆动，意味着系数将远离 0 (参见图 4-8)；当加上这个罚时，我们强制更平坦、更平滑的拟合。罚多少依赖于 λ ，它是先验方差的逆，即我们期望先验的权重离 0 多远。也就是说，有这样的先验等价于迫使参数接近于 0。我们将在第 14 章更详细地讨论这一问题。

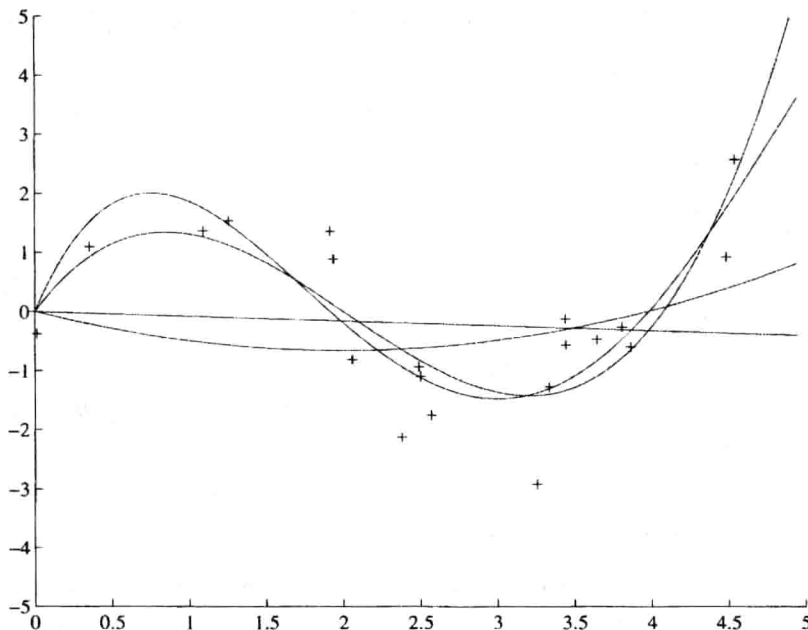


图 4-8 图 4-8 与图 4-5 同样的设置，拟合 1~4 阶多项式。系数的大小随多项式阶数的增加而增加，它们是 1: $[-0.0769, 0.00116]^T$, 2: $[0.1682, -0.6657, 0.0080]^T$, 3: $[0.4238, -2.5778, 3.4675, -0.0002]^T$, 4: $[-0.1093, 1.4356, -5.5007, 6.0454, -0.0019]^T$

当选择先验使得最简单的模型具有最高概率时(根据奥卡姆剃刀规则),贝叶斯方法、正则化、SRM 和 MDL 是等价的。交叉验证与其他模型选择的方法不同,因为它不对模型做任何先验假设。如果有足够大的验证数据集,它就是最好的方法。在数据样本很小时,其他模型变得有用。

4.9 注释

最大似然和贝叶斯估计基础的一个好资源是 Ross 1987。许多模式识别教材都讨论参数模型分类(例如,MacLachlan 1992; Devroye、Györfi 和 Lugosi 1996; Webb 1999; Duda、Hart 和 Stork 2001)。检查一元正态性的检验可以在 Rencher 1995 中找到。

Geman、Bienenstock 和 Doursat (1992) 讨论了几个学习模型的偏倚和方差分解,这些我们也将后边的章节讨论。偏倚/方差分解是针对回归的;对于 0/1 损失,误差的这种偏倚、方差和噪声的加法分解是不可能的,因为对于分类,如果意外落到边界的另一侧就会出错。对于两类问题,如果正确的后验是 0.7,而估计是 0.8,则没有错误,仅当估计小于 0.5 时才出错。对于分类,不同的研究者提出了偏倚和方差的不同定义,综述参见 Friedman 1997。

4.10 习题

1. 写出产生以给定的 p 为参数的伯努利样本的程序,并且写出由样本计算 \hat{p}_i 的程序。
2. 写出多项样本的对数似然,并证明式(4.6)。
3. 写出产生以给定 μ, σ 为参数的正态样本并由样本计算 m 和 s 的程序。对 μ 假定先验分布,用贝叶斯估计做相同的工作。
4. 给定两个正态分布 $p(x | C_1) \sim \mathcal{N}(\mu_1, \sigma_1^2)$ 和 $p(x | C_2) \sim \mathcal{N}(\mu_2, \sigma_2^2)$ 以及 $P(C_1)$ 和 $P(C_2)$, 解析地计算贝叶斯判别点。
5. 在高斯密度的情况下,似然比 $\frac{p(x | C_1)}{p(x | C_2)}$ 是什么?
6. 对于两类问题,用不同的方差为两个类产生正态样本,然后使用参数分类法估计判别点。将它与理论值进行比较。
7. 假定一个线性模型,然后加入 0 均值的高斯噪声来产生一个样本。把样本一分为二,分别作为训练集和验证集。在训练集的这一半上使用线性回归。在验证集上计算误差。对 2 次和 3 次多项式进行同样的处理。
8. 当训练集较小时,方差对误差的贡献可能比偏倚大。在这种情况下,我们可能更喜欢简单模型,即使我们知道对于我们的任务它太简单。你能给出一个例子吗?
9. 假设给定样本 $X_i = \{x_i', r_i'\}$, 我们定义 $g_i(x) = r_i'$, 即我们对任意 x 的估计是数据集 X_i 的第一个实例(未排序)的 r 值。与 $g_i(x) = 2$ 和 $g_i(x) = \sum_i r_i' / N$ 相比,关于它的偏倚和方差你有何评论? 如果样本是有序的并使 $g_i(x) = \min_i r_i'$, 情况又如何?
10. 在式(4.40)中,改变 λ 对偏倚和方差的影响是什么?

83
1
84

82

4.11 参考文献

- Andrieu, C., N. de Freitas, A. Doucet, and M. I. Jordan. 2003. "An Introduction to MCMC for Machine Learning." *Machine Learning* 50: 5-43.
- Breiman, L. 1998. "Bias-Variance, Regularization, Instability and Stabilization." In *Neural Networks and Machine Learning*, ed. C. M. Bishop, 27-56. Berlin: Springer.
- Devroye, L., L. Györfi, and G. Lugosi. 1996. *A Probabilistic Theory of Pattern Recognition*. New York: Springer.
- Duda, R. O., P. E. Hart, and D. G. Stork. 2001. *Pattern Classification*, 2nd ed. New York: Wiley.
- Friedman, J. H. 1997. "On Bias, Variance, 0/1-Loss and the Curse of Dimensionality." *Data Mining and Knowledge Discovery* 1: 55-77.
- Geman, S., E. Bienenstock, and R. Doursat. 1992. "Neural Networks and the Bias/Variance Dilemma." *Neural Computation* 4: 1-58.
- Grünwald, P. D. 2007. *The Minimum Description Length Principle*. Cambridge, MA: MIT Press.
- Hastie, T., R. Tibshirani, and J. Friedman. 2001. *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. New York: Springer.
- McLachlan, G. J. 1992. *Discriminant Analysis and Statistical Pattern Recognition*. New York: Wiley.
- Rencher, A. C. 1995. *Methods of Multivariate Analysis*. New York: Wiley.
- Rissanen, J. 1978. "Modeling by Shortest Data Description." *Automatica* 14: 465-471.
- Ross, S. M. 1987. *Introduction to Probability and Statistics for Engineers and Scientists*. New York: Wiley.
- Vapnik, V. 1995. *The Nature of Statistical Learning Theory*. New York: Springer.
- Webb, A. 1999. *Statistical Pattern Recognition*. London: Arnold.

多元方法

在第4章,我们讨论了分类和回归的参数方法。现在,我们将它们推广到多元情况,其中我们有多个输入,并且输出(即类编码或连续输出)是这些输入的函数。这些输入可能是离散的或数值的。我们将讨论如何从标记的多元样本学习这样的函数,以及如何根据已有数据调整学习方法的复杂度。

5.1 多元数据

在许多应用中,许多测量都在每个个体或者事件上进行,并产生一个观测向量。样本可以看作一个数据矩阵(data matrix)

$$\mathbf{X} = \begin{bmatrix} X_1^1 & X_2^1 & \cdots & X_d^1 \\ X_1^2 & X_2^2 & \cdots & X_d^2 \\ \vdots & \vdots & \ddots & \vdots \\ X_1^N & X_2^N & \cdots & X_d^N \end{bmatrix}$$

其中 d 列对应 d 个变量,表示在个体或事件上的测量结果。它们也称为输入(input)、特征(feature)或属性(attribute)。 N 行对应在 N 个个体或事件上的独立同分布的观测(observation)、样例(example)或实例(instance)。

例如,在对贷款申请做决定时,观测向量是与客户相关的一些信息,包括客户的年龄、婚姻状况、年收入等,并且我们有 N 个这样的老用户。这些测量也许有不同的尺度,例如,年龄用年计算,年收入用货币单位计算。某些(如年龄)可能是数值的,某些(如婚姻状况)可能是离散的。

通常,这些变量是相关的。如果它们不相关,就没有必要做多元分析。我们的目标也许是化简(simplification),也就是用相对少的参数汇总大量数据。我们的目标也许是探测性的(exploratory),并且我们可能对产生关于数据的假设感兴趣。在有些应用中,我们对由其他变量的值预测一个变量值很感兴趣。如果被预测变量是离散的,这就是多元分类问题;如果是数值的,这就是多元回归问题。

5.2 参数估计

均值向量(mean vector) $\boldsymbol{\mu}$ 的每个元素都是 \mathbf{X} 一个列的均值:

$$E[\mathbf{x}] = \boldsymbol{\mu} = [\mu_1, \dots, \mu_d]^T \quad (5.1)$$

X_i 的方差记作 σ_i^2 , 两个变量 X_i 和 X_j 的协方差定义为

$$\sigma_{ij} \equiv \text{Cov}(X_i, X_j) = E[(X_i - \mu_i)(X_j - \mu_j)] = E[X_i X_j] - \mu_i \mu_j \quad (5.2)$$

满足 $\sigma_{ij} = \sigma_{ji}$, 并且当 $i = j$ 时, $\sigma_{ii} = \sigma_i^2$ 。 d 个变量就有 d 个方差和 $d(d-1)/2$ 个协方差。通

常表示为 $d \times d$ 矩阵, 称为协方差矩阵 (covariance matrix), 用 Σ 表示, 其第 (i, j) 个元素是 σ_{ij} :

$$\Sigma = \begin{bmatrix} \sigma_1^2 & \sigma_{12} & \cdots & \sigma_{1d} \\ \sigma_{21} & \sigma_2^2 & \cdots & \sigma_{2d} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{d1} & \sigma_{d2} & \cdots & \sigma_d^2 \end{bmatrix}$$

对角线上的元素是方差, 非对角线上的元素是协方差, 并且矩阵是对称的。使用向量矩阵记号

$$\Sigma \equiv \text{Cov}(\mathbf{X}) = E[(\mathbf{X} - \boldsymbol{\mu})(\mathbf{X} - \boldsymbol{\mu})^T] = E[\mathbf{X}\mathbf{X}^T] - \boldsymbol{\mu}\boldsymbol{\mu}^T \quad (5.3)$$

如果两个向量是线性相关的, 则协方差是正还是负取决于线性关系的斜率是正还是负。但是相关性的大小很难解释, 因为它取决于两个变量的测量单位。变量 X_i 和 X_j 的相关性 (correlation) 是一个规范化到 -1 到 +1 之间的统计量, 定义为

$$\text{Corr}(X_i, X_j) \equiv \rho_{ij} = \frac{\sigma_{ij}}{\sigma_i \sigma_j} \quad (5.4)$$

如果两个变量是相互独立的, 那么其协方差为 0, 因而相关性为 0。然而, 其逆不真: 变量也许是相关的 (以非线性方式), 并且它们的相关性可能为 0。

给定多元样本, 可以计算这些参数的估计: 均值的最大似然估计是样本均值 \mathbf{m} 。它的第 i 维是 \mathbf{X} 的第 i 列的平均值:

$$\mathbf{m} = \frac{\sum_{t=1}^N \mathbf{x}^t}{N}, \quad \text{其中 } m_i = \frac{\sum_{t=1}^N x_i^t}{N}, \quad i = 1, \dots, d \quad (5.5)$$

Σ 的估计是样本协方差 (sample covariance) 矩阵 S , 其元素是

$$s_i^2 = \frac{\sum_{t=1}^N (x_i^t - m_i)^2}{N} \quad (5.6)$$

$$s_{ij} = \frac{\sum_{t=1}^N (x_i^t - m_i)(x_j^t - m_j)}{N} \quad (5.7)$$

它们是有偏估计, 但如果在应用中估计的变化显著依赖于我们用 N 还是用 $N-1$ 来除, 那么我们将遇到严重的麻烦。

样本相关 (sample correlation) 系数

$$r_{ij} = \frac{s_{ij}}{s_i s_j} \quad (5.8)$$

而样本的相关矩阵 \mathbf{R} 包含 r_{ij} 。

5.3 缺失值估计

观测中的某些变量的值可能缺失。最好的策略是把这些观测一同丢弃, 但是, 一般我们没有足够大的样本来让我们这样做, 并且我们不想丢弃数据, 因为非缺失的条目确实包含信息。我们试图通过估计它们来填写缺失的条目, 这被称作估算 (imputation)。

在均值估算 (mean imputation) 中, 对于数值变量, 我们用现有数据的均值 (平均值) 来代替样本中的变量值。对离散变量, 我们用最可能出现的值, 即数据中出现最多的值, 来填写

变量值。

在回归估算(imputation by regression)中,我们试图从其他已知变量的值来预测缺失的变量值。根据缺失变量的类型,我们分别定义一个回归或分类问题,用其值已知的数据点训练。如果许多不同变量都缺失,则我们取均值作为初始估计,并且反复执行该过程直到被预测的值稳定。如果这些变量不是高度相关的,则回归方法与均值估算等价。

然而,根据环境,有时特定属性值的缺失也许很重要。例如,在信用卡申请中,如果申请人不提供他或她的电话号码,那也许是一条关键信息。在这样的情况下,我们用一个单独的值表示它,指明该值缺失并照此使用。

5.4 多元正态分布

在多元情况下,其中 \mathbf{x} 是 d 维、正态分布的,我们有

$$p(\mathbf{x}) = \frac{1}{(2\pi)^{d/2} |\mathbf{\Sigma}|^{1/2}} \exp \left[-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \mathbf{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right] \quad (5.9)$$

并且我们记 $\mathbf{x} \sim \mathcal{N}_d(\boldsymbol{\mu}, \mathbf{\Sigma})$,其中 $\boldsymbol{\mu}$ 是均值向量, $\mathbf{\Sigma}$ 是协方差矩阵(参见图5-1)。正如

$$\frac{(\mathbf{x} - \boldsymbol{\mu})^2}{\sigma^2} = (\mathbf{x} - \boldsymbol{\mu}) (\sigma^2)^{-1} (\mathbf{x} - \boldsymbol{\mu})$$

是 \mathbf{x} 到 $\boldsymbol{\mu}$ 的以标准差为单位、对不同的方差规范化的平方距离一样,在多元情况下,使用马氏距离(Mahalanobis distance):

$$(\mathbf{x} - \boldsymbol{\mu})^T \mathbf{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \quad (5.10)$$

$(\mathbf{x} - \boldsymbol{\mu})^T \mathbf{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) = c^2$ 是以 $\boldsymbol{\mu}$ 为中心的 d 维超椭球,并且它的形状和方向由 $\mathbf{\Sigma}$ 决定。

90 由于使用了 $\mathbf{\Sigma}$ 的逆,所以如果一个变量比其他变量的方差大,则它在马氏距离中的权重较小。类似地,两个高度相关变量的贡献没有两个相关性较低变量的贡献大。这样,使用协方差矩阵的逆具有将所有变量标准化(具有单位方差)并消除相关性的效果。

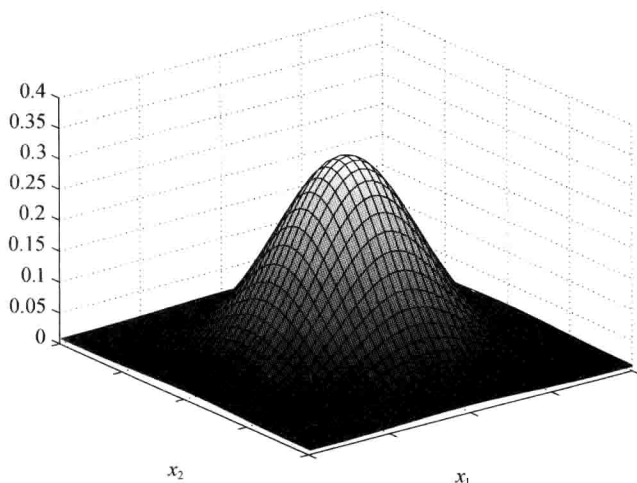


图 5-1 二元正态分布

为便于显示,让我们考虑二元情况,其中 $d=2$ (参见图5-2)。当变量独立时,密度的主轴与输入轴平行。如果方差不同,则密度变成椭圆。密度根据协方差(相关性)的符号而旋

转。均值向量为 $\boldsymbol{\mu}^T = [\mu_1, \mu_2]$ ，协方差矩阵通常表示为

$$\boldsymbol{\Sigma} = \begin{bmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 \end{bmatrix}$$

$\text{Cov}(x_1, x_2) = 0, \text{Var}(x_1) = \text{Var}(x_2)$ $\text{Cov}(x_1, x_2) = 0, \text{Var}(x_1) > \text{Var}(x_2)$

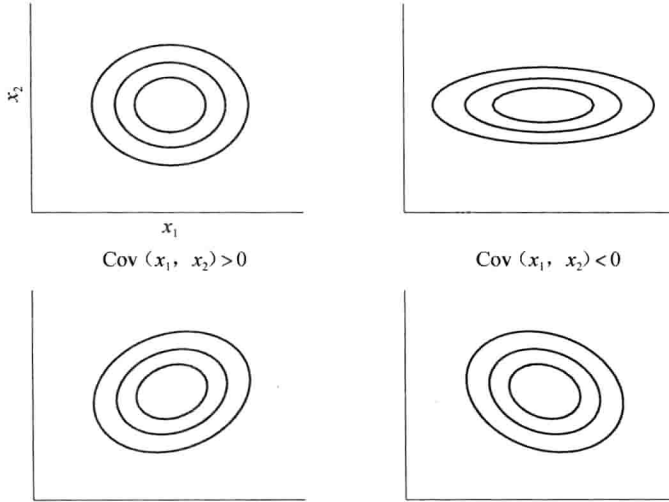


图 5-2 二元正态分布的等概率围线图。其中心由均值给定，其形状和方向依赖于协方差矩阵

二元联合密度可以表示为如下形式(见习题 1)

$$p(x_1, x_2) = \frac{1}{2\pi\sigma_1\sigma_2\sqrt{1-\rho^2}} \exp\left[-\frac{1}{2(1-\rho^2)}(z_1^2 - 2\rho z_1 z_2 + z_2^2)\right] \quad (5.11)$$

其中, $z_i = (x_i - \mu_i)/\sigma_i (i = 1, 2)$ 是标准化变量, 称为 z -规范化 (z -normalization)。记住, 当 $|\rho| < 1$ 时,

$$z_1^2 + 2\rho z_1 z_2 + z_2^2 = \text{常数}$$

是椭圆方程。当 $\rho > 0$ 时, 椭圆的主轴有正斜率, 当 $\rho < 0$ 时, 主轴有负斜率。

在式(5.11)的扩展马氏距离中, 每个变量都被规范化, 使具有单位方差, 交叉项修正了两个变量之间的相关性。

概率密度依赖于五个参数: 两个均值、两个方差和相关性。只要方差非 0 并且 $|\boldsymbol{\Sigma}| < 1$, 则就是非奇异的, 因此是正定的。如果 ρ 是 +1 或者 -1, 则两个变量是线性相关的, 观测事实上是一维的, 两个变量中的一个可以去掉。如果 $\rho = 0$, 则两个变量是独立的, 交叉项消失, 我们得到两个一元密度的积。

在多元情况下, 小的 $|\boldsymbol{\Sigma}|$ 值表明样本接近 $\boldsymbol{\mu}$, 正如单变量情况, 小的 σ^2 表示样本接近 μ 。小的 $|\boldsymbol{\Sigma}|$ 还表示两个变量之间的高度相关性。 $\boldsymbol{\Sigma}$ 是对称正定矩阵; 这是 $\text{Var}(X) > 0$ 的多元说法。否则, $\boldsymbol{\Sigma}$ 是奇异的, 它的行列式的值为 0。这要么是由于维之间的线性依赖性, 要么是因为有一维具有 0 方差。在这种情况下, 应该将维度归约为正定矩阵; 第 6 章将讨论这个问题的处理方法。

如果 $\mathbf{x} \sim \mathcal{N}_d(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, 则 \mathbf{x} 的每个维都是一元正态的。(其逆不真: 每一个 X_i 都可以是一元正态的, 而 \mathbf{X} 不一定是多元正态的。)实际上, 变量的任意 $k < d$ 的子集都是 k 元正态的。

一个特殊的情况是, \mathbf{x} 的分量是独立的, 并且 $i \neq j$ 时 $\text{Cov}(X_i, X_j) = 0$, 并且 $\text{Var}(X_i) = \sigma_i^2$, $\forall i$. 于是, 协方差矩阵是对角的, 联合密度是各一元密度的乘积:

$$p(\mathbf{x}) = \prod_{i=1}^d p_i(x_i) = \frac{1}{(2\pi)^{d/2} \prod_{i=1}^d \sigma_i} \exp \left[-\frac{1}{2} \sum_{i=1}^d \left(\frac{x_i - \mu_i}{\sigma_i} \right)^2 \right] \quad (5.12)$$

现在, 我们考察另一个性质, 它将在以后的章节中用到。我们假设 $\mathbf{x} \sim \mathcal{N}_d(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, 其中 $\mathbf{w} \in \mathbb{R}^d$, 则

$$\mathbf{w}^T \mathbf{x} = w_1 x_1 + w_2 x_2 + \cdots + w_d x_d \sim \mathcal{N}(\mathbf{w}^T \boldsymbol{\mu}, \mathbf{w}^T \boldsymbol{\Sigma} \mathbf{w})$$

给定

$$\begin{aligned} E[\mathbf{w}^T \mathbf{x}] &= \mathbf{w}^T E[\mathbf{x}] = \mathbf{w}^T \boldsymbol{\mu} \\ \text{Var}(\mathbf{w}^T \mathbf{x}) &= E[(\mathbf{w}^T \mathbf{x} - \mathbf{w}^T \boldsymbol{\mu})^2] = E[(\mathbf{w}^T \mathbf{x} - \mathbf{w}^T \boldsymbol{\mu})(\mathbf{w}^T \mathbf{x} - \mathbf{w}^T \boldsymbol{\mu})] \\ &= E[\mathbf{w}^T (\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T \mathbf{w}] = \mathbf{w}^T E[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T] \mathbf{w} \\ &= \mathbf{w}^T \boldsymbol{\Sigma} \mathbf{w} \end{aligned} \quad (5.13)$$

这就是说, d 维正态分布在向量 \mathbf{w} 上的投影是一元正态分布。在一般情况下, 如果 \mathbf{W} 是 $d \times k$ 矩阵, 其秩 $k < d$, 则 k 维 $\mathbf{W}^T \mathbf{x}$ 矩阵是 k 元正态分布:

$$\mathbf{W}^T \mathbf{x} \sim \mathcal{N}_k(\mathbf{W}^T \boldsymbol{\mu}, \mathbf{W}^T \boldsymbol{\Sigma} \mathbf{W}) \quad (5.14)$$

也就是说, 如果我们把一个 d 维正态分布投影到 k 维的空间, 则投影是 k 维正态分布。

93

5.5 多元分类

当 $\mathbf{x} \in \mathbb{R}^d$ 时, 如果取类条件密度 $p(\mathbf{x} | C_i)$ 为正态密度 $\mathcal{N}_d(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$, 则我们有

$$p(\mathbf{x} | C_i) = \frac{1}{(2\pi)^{d/2} |\boldsymbol{\Sigma}_i|^{1/2}} \exp \left[-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_i)^T \boldsymbol{\Sigma}_i^{-1} (\mathbf{x} - \boldsymbol{\mu}_i) \right] \quad (5.15)$$

这样做的主要原因是它分析简单(Duda、Hart 和 Stork 2001)。此外, 正态分布密度函数是许多自然现象的模型, 因为大多数类的样本都可以看作是简单原型 $\boldsymbol{\mu}_i$ 的轻微改变版本, 并且协方差矩阵 $\boldsymbol{\Sigma}_i$ 表示每个变量中的噪声量和这些噪声源的相关性。尽管真实数据可能常常并非严格多元正态的, 但是这是一个有用的近似。除了它易于进行数学处理外, 该模型对偏离正态分布的鲁棒性在许多工作中都展示出来(例如 McLachlan 1992)。然而, 一个明显的要求是一个类的样本应该形成单个组; 如果有多个组, 应该使用混合模型(第7章)。

假设我们要预测顾客可能感兴趣的汽车类型。不同的汽车是类, 而 \mathbf{x} 是顾客的可观测数据, 例如年龄和收入。 $\boldsymbol{\mu}_i$ 是购买 i 类汽车的顾客年龄和收入的均值向量, 而 $\boldsymbol{\Sigma}_i$ 是它们的协方差矩阵: σ_{i1}^2 和 σ_{i2}^2 分别是年龄和收入的方差, 并且 σ_{i12} 是购买 i 类汽车的顾客年龄和收入的协方差。

当我们定义判别式函数为

$$g_i(\mathbf{x}) = \log p(\mathbf{x} | C_i) + \log P(C_i)$$

并假定 $p(\mathbf{x} | C_i) \sim \mathcal{N}_d(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$ 时, 我们有

$$g_i(\mathbf{x}) = -\frac{d}{2} \log 2\pi - \frac{1}{2} \log |\boldsymbol{\Sigma}_i| - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_i)^T \boldsymbol{\Sigma}_i^{-1} (\mathbf{x} - \boldsymbol{\mu}_i) + \log P(C_i) \quad (5.16)$$

给定 $K \geq 2$ 个类的训练样本 $\mathcal{X} = \{\mathbf{x}^t, \mathbf{r}^t\}$, 其中如果 $\mathbf{x}^t \in C_i$, 则 $r_i^t = 1$, 否则为 0。分别对每个类求最大似然, 从而找到均值和协方差的估计:

$$\hat{P}(C_i) = \frac{\sum_t r_i^t}{N} \quad (5.18)$$

$$\mathbf{m}_i = \frac{\sum_t r_i^t \mathbf{x}^t}{\sum_t r_i^t}$$

$$\mathbf{S}_i = \frac{\sum_t r_i^t (\mathbf{x}^t - \mathbf{m}_i)(\mathbf{x}^t - \mathbf{m}_i)^T}{\sum_t r_i^t}$$

94

然后, 将这些代入判别式函数, 得到判别式的估计。忽略第一个常数项, 我们有

$$g_i(\mathbf{x}) = -\frac{1}{2} \log |\mathbf{S}_i| - \frac{1}{2} (\mathbf{x} - \mathbf{m}_i)^T \mathbf{S}_i^{-1} (\mathbf{x} - \mathbf{m}_i) + \log \hat{P}(C_i) \quad (5.19)$$

把它展开, 我们得到

$$g_i(\mathbf{x}) = -\frac{1}{2} \log |\mathbf{S}_i| - \frac{1}{2} (\mathbf{x}^T \mathbf{S}_i^{-1} \mathbf{x} - 2 \mathbf{x}^T \mathbf{S}_i^{-1} \mathbf{m}_i + \mathbf{m}_i^T \mathbf{S}_i^{-1} \mathbf{m}_i) + \log \hat{P}(C_i)$$

它定义了一个二次判别式(quadratic discriminant)(参见图 5-3), 也可以写作

$$g_i(\mathbf{x}) = \mathbf{x}^T \mathbf{W}_i \mathbf{x} + \mathbf{w}_i^T \mathbf{x} + w_{i0} \quad (5.20)$$

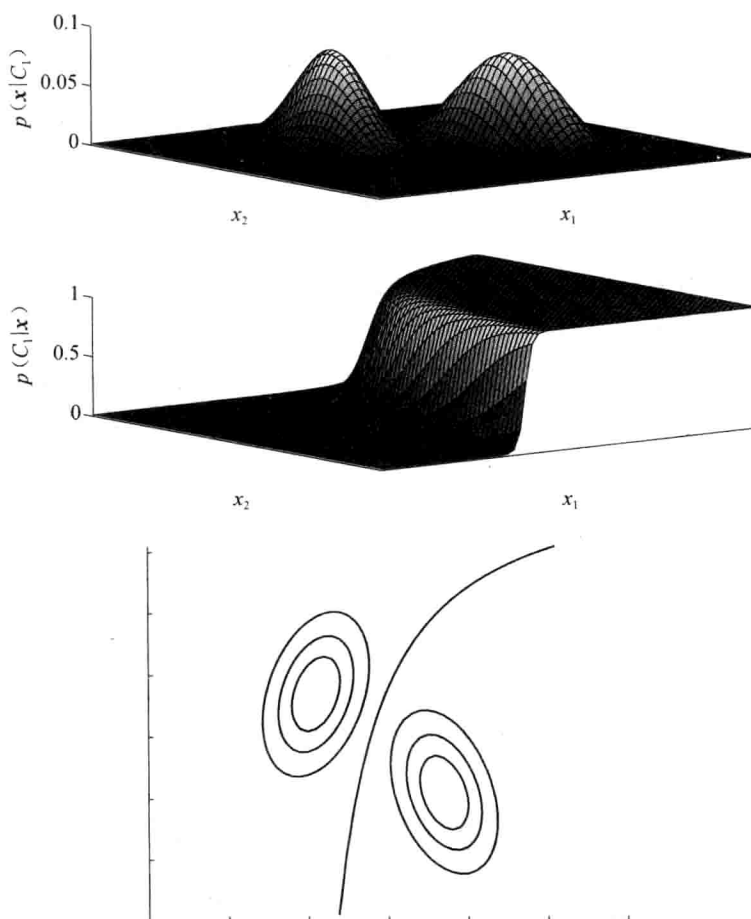


图 5-3 类具有不同的协方差矩阵。其中一个是类的似然密度和后验概率(顶部)。类分布用等概率围线表示, 并且绘出判别式(底部)

其中

$$\begin{aligned} \mathbf{W}_i &= -\frac{1}{2}\mathbf{S}_i^{-1} \\ \mathbf{w}_i &= \mathbf{S}_i^{-1}\mathbf{m}_i \\ w_{i0} &= -\frac{1}{2}\mathbf{m}_i^T\mathbf{S}_i^{-1}\mathbf{m}_i - \frac{1}{2}\log|\mathbf{S}_i| + \log\hat{P}(C_i) \end{aligned}$$

对于均值, 需要估计的参数数量为 $K \cdot d$ 个, 而对于协方差矩阵为 $K \cdot d(d+1)/2$ 个。当 d 大而样本小时, \mathbf{S}_i 可能是奇异的, 并且其逆可能不存在。或者, $|\mathbf{S}_i|$ 可能非零但是太小, 这种情况会不稳定; \mathbf{S}_i 的小变化会引起 \mathbf{S}_i^{-1} 的大变化。为了使小样本上的估计可靠, 我们可能希望通过重新设计特性提取算法并选择特征子集, 或者组合已有特征来降低维度 d 。我们将在第6章讨论这样的方法。

另外一个可能的做法是汇集数据, 并且对所有的类估计公共协方差矩阵:

$$\mathbf{S} = \sum_i \hat{P}(C_i)\mathbf{S}_i \quad (5.21)$$

在相同协方差矩阵的情况下, 式(5.19)化简为

$$g_i(\mathbf{x}) = -\frac{1}{2}(\mathbf{x} - \mathbf{m}_i)^T\mathbf{S}^{-1}(\mathbf{x} - \mathbf{m}_i) + \log\hat{P}(C_i) \quad (5.22)$$

对于均值, 参数数量为 $K \cdot d$ 个, 而对于共享协方差矩阵为 $d(d+1)/2$ 个。如果先验相等, 则最佳决策规则是把输入指派到与输入均值的马氏距离最小的类。像前面一样, 不相等的先验将边界移向不太可能的类。注意在这种情况下, 二次项 $\mathbf{x}^T\mathbf{S}^{-1}\mathbf{x}$ 被取消, 因为它出现在所有的判别式中, 并且决策边界是线性的, 导致如下线性判别式(linear discriminant)(图5-4):

$$g_i(\mathbf{x}) = \mathbf{w}_i^T\mathbf{x} + w_{i0} \quad (5.23)$$

其中

$$\begin{aligned} \mathbf{w}_i &= \mathbf{S}^{-1}\mathbf{m}_i \\ w_{i0} &= -\frac{1}{2}\mathbf{m}_i^T\mathbf{S}^{-1}\mathbf{m}_i + \log\hat{P}(C_i) \end{aligned}$$

这种线性分类器的决策区域是凸的, 即在一个决策区域内任意选择两个点并用一条直线连接, 直线上的所有点都落在该区域内。

通过假定协方差矩阵的所有非对角线元素均为零, 从而假定变量都是独立的, 问题可以进一步简化。这是朴素贝叶斯分类(naive Bayes' classifier), 其中 $p(x_j | C_i)$ 是一元高斯的。 \mathbf{S} 和它的逆都是对角的, 并且我们有

$$g_i(\mathbf{x}) = -\frac{1}{2}\sum_{j=1}^d \left(\frac{x_j^i - m_{ij}}{s_j} \right)^2 + \log\hat{P}(C_i) \quad (5.24)$$

项 $((x_j^i - m_{ij})/s_j)^2$ 有规范化作用并以标准差为单位度量距离。从几何学角度来说, 类是超椭圆体, 并且因为协方差为零, 它还是轴对齐的(见图5-5)。对于均值, 参数的数量为 $K \cdot d$, 而对于方差为 d 。这样, \mathbf{S} 的复杂度由 $\mathcal{O}(d^2)$ 降低为 $\mathcal{O}(d)$ 。

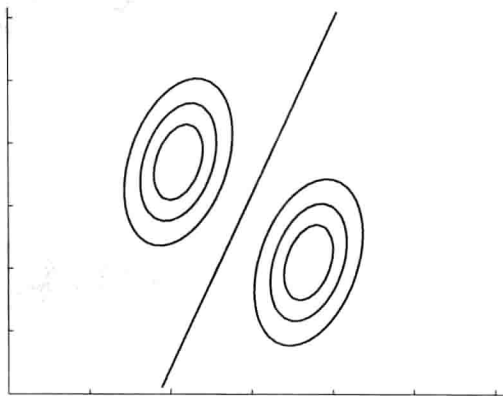


图5-4 协方差可以是任意的, 但是被两个类共享

如果我们假定所有的变量是等同的, 马氏距离归结为欧氏距离 (Euclidean distance), 问题还可以进一步简化。在几何图像上, 分布是球形的, 并以均值向量 \mathbf{m}_i 为中心 (见图 5-6)。于是, $|\mathbf{S}| = s^{2d}$ 并且 $\mathbf{S}^{-1} = (1/s^2)\mathbf{I}$ 。参数是 $K \cdot d$ 个均值和一个 s^2 。

$$g_i(\mathbf{x}) = -\frac{\|\mathbf{x} - \mathbf{m}_i\|^2}{2s^2} + \log \hat{P}(C_i) = -\frac{1}{2s^2} \sum_{j=1}^d (x_j' - m_{ij})^2 + \log \hat{P}(C_i) \quad (5.25)$$

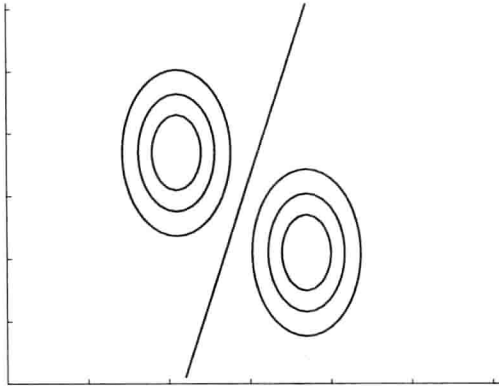


图 5-5 所有的类都具有相等的对角协方差, 但是方差不相等

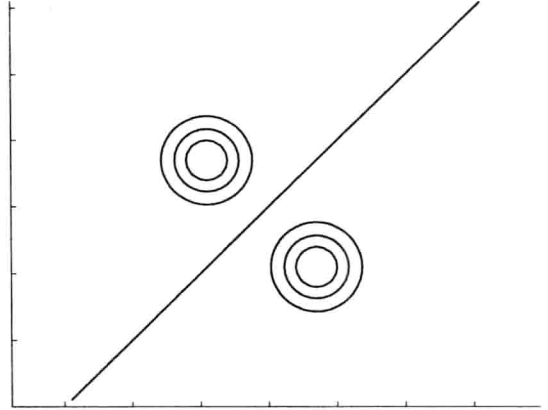


图 5-6 所有类具有相等的、在两个维上均具有相等方差的对角协方差矩阵

如果先验相等, 则我们有 $g_i(\mathbf{x}) = -\|\mathbf{x} - \mathbf{m}_i\|^2$ 。这被称为最近均值分类 (nearest mean classifier), 因为它把输入指派到最近均值的类。如果每个均值都被看作是类的理想原型或模板, 那么这就是模板匹配 (template matching) 过程。它可以被扩展为

$$\begin{aligned} g_i(\mathbf{x}) &= -\|\mathbf{x} - \mathbf{m}_i\|^2 = -(\mathbf{x} - \mathbf{m}_i)^T(\mathbf{x} - \mathbf{m}_i) \\ &= -(\mathbf{x}^T\mathbf{x} - 2\mathbf{m}_i^T\mathbf{x} + \mathbf{m}_i^T\mathbf{m}_i) \end{aligned} \quad (5.26)$$

第一项 $\mathbf{x}^T\mathbf{x}$ 出现在所有的 $g_i(\mathbf{x})$ 中, 可以去掉, 并且我们可以把判别式函数写成

$$g_i(\mathbf{x}) = \mathbf{w}_i^T\mathbf{x} + w_{i0} \quad (5.27)$$

其中 $\mathbf{w}_i = \mathbf{m}_i$, $w_{i0} = -(1/2)\|\mathbf{m}_i\|^2$ 。如果所有的 \mathbf{m}_i 有相似的范数, 则 w_{i0} 也可以忽略, 并且我们可以使用

$$g_i(\mathbf{x}) = \mathbf{m}_i^T\mathbf{x} \quad (5.28)$$

当 \mathbf{m}_i 的范数可比较时, 也可以使用点积代替 (负的) 欧氏距离作为相似性度量。

我们实际上可以把寻找最佳判别函数的任务看作是寻找最佳距离函数。这可以被看作是另外一种分类方法: 我们不是要学习判别式函数 $g_i(\mathbf{x})$, 而是要学习一个合适的距离函数 $\mathcal{D}(\mathbf{x}_1, \mathbf{x}_2)$, 使得对任意 $\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3$, 其中 $\mathbf{x}_1, \mathbf{x}_2$ 属于相同的类, 而 $\mathbf{x}_1, \mathbf{x}_3$ 属于两个不同的类, 我们希望能有

$$\mathcal{D}(\mathbf{x}_1, \mathbf{x}_2) < \mathcal{D}(\mathbf{x}_1, \mathbf{x}_3)$$

5.6 调整复杂度

在表 5-1 中, 我们看到如何减少协方差矩阵的参数数目, 在简单模型的适用性和通用性之间折中。这是偏倚/方差两难选择的又一个例子。当我们做简化协方差矩阵的假设并降低

被估计的参数数目时,我们就有引进偏倚的风险(参见图 5-7)。另一方面,如果不做这种假设,并且矩阵是任意的,则二次判别式函数在小数据集上会有很大的方差。理想情况取决于已有数据所表示的问题的复杂度和我们所拥有的数据的规模。当我们拥有小数据集时,尽管协方差矩阵不同,但是假定共享协方差矩阵也许更好;单个协方差矩阵具有较少参数,并且可以利用更多的数据来估计,即用所有类的实例估计。这相当于使用线性判别式(linear discriminant)。分类经常用到线性判别式,我们将在第 10 章更详细地讨论它。

表 5-1 通过简化假设降低方差

假 设	协方差矩阵	参数数目
共享、超球	$S_i = S = s^2 I$	1
共享、轴对齐	$S_i = S$, 其中 $s_{ij} = 0$	d
共享、超椭圆	$S_i = S$	$d(d+1)/2$
不同、超椭圆	S_i	$K \cdot (d(d+1)/2)$

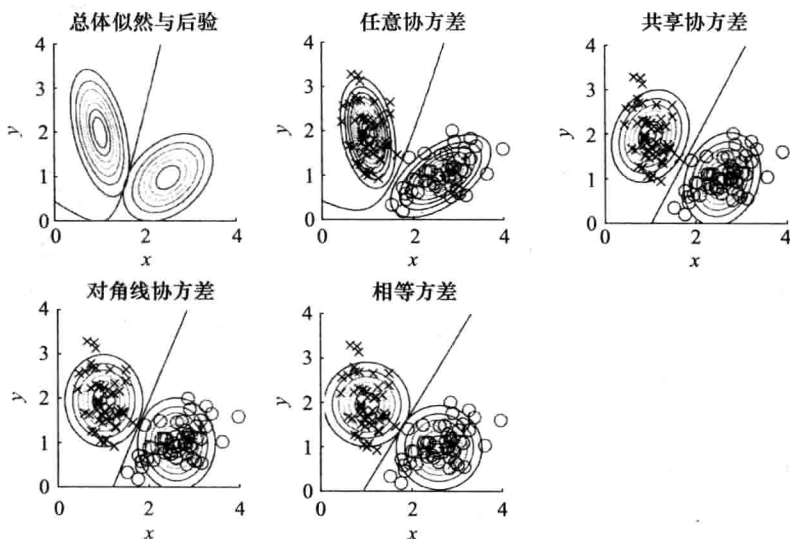


图 5-7 拟合相同数据的不同协方差矩阵导致不同的边界

注意,当我们用欧氏距离度量相似性时,我们假设所有的变量都具有相同的方差,并且它们是相互独立的。在许多情况下,这并不成立。例如,年龄与年收入具有不同的单位,并且在许多情况下是相关的。在这种情况下,可以在预处理阶段对输入进行 z-规范化(具有 0 均值和单位方差),然后使用欧氏距离。另一方面,有些时候即使变量是相关的,如果我们没有足够的准确数据计算相关性,也许最好还是假设它们不相关,并使用朴素贝叶斯分类。

Friedman (1989) 提出一个组合所有这些特殊情况的方法,称为正则化判别式分析(regularized discriminant analysis, RDA)。我们还记得,正则化方法对应于从大方差开始到小方差的方法,有增加偏倚的风险。在利用高斯密度的参数化分类情况下,协方差矩阵可以表示成三种特殊情况的加权平均:

$$S'_i = \alpha \sigma^2 I + \beta S + (1 - \alpha - \beta) S_i \quad (5.29)$$

当 $\alpha = \beta = 0$ 时,我们得到二次分类器。当 $\alpha = 0, \beta = 1$ 时,协方差矩阵被共享,我们得到线性分类器。当 $\alpha = 1, \beta = 0$ 时,协方差矩阵是对角阵, σ^2 在对角线上,我们得到最

近均值分类。在这些极端情况之间，我们得到所有的不同分类方法，其中 α 和 β 通过交叉确认优化。

当数据集较小时，另一种正则化方法是通过定义 μ_i 和 S_i 上的先验，使用贝叶斯方法，或者使用交叉验证选择表 5-1 中给出的四种情况中最好者。

5.7 离散特征

在许多应用中，我们有取 n 个不同值的离散属性。例如，一个属性可能是颜色 $\in \{\text{红, 蓝, 绿, 黑}\}$ ，另外一个可能是像素 $\in \{\text{有, 无}\}$ 。我们假设 x_j 是二元的(伯努利)，其中

$$p_{ij} \equiv p(x_j = 1 | C_i)$$

如果 x_j 是独立的二元变量，则我们有

$$p(\mathbf{x} | C_i) = \prod_{j=1}^d p_{ij}^{x_j} (1 - p_{ij})^{(1-x_j)}$$

这是朴素贝叶斯分类的另一个例子，其中 $p(x_j | C_i)$ 是伯努利分布。判别式函数是

$$\begin{aligned} g_i(\mathbf{x}) &= \log p(\mathbf{x} | C_i) + \log P(C_i) \\ &= \sum_j [x_j \log p_{ij} + (1 - x_j) \log(1 - p_{ij})] + \log P(C_i) \end{aligned} \quad (5.30)$$

它是线性的。 p_{ij} 的估计是

$$\hat{p}_{ij} = \frac{\sum_t x_j^t r_i^t}{\sum_t r_i^t} \quad (5.31)$$

这种方法用于文档分类(document categorization)。文档分类的一个例子是把新闻报道分成不同的类别，如政治、体育、时尚等类别。在词袋(bag of words)表示中，预先选择 d 个我们相信能够提供与类相关的信息的词(Manning 和 Schütze 1999)。例如，在新闻分类中，有用的词是诸如 missile、athlete 和 couture 这样的词，而不是像 model 甚至是 runaway 这样有歧义的词。在这种表示下，每个文本是一个 d 维二元向量，其中如果第 j 个词出现在该文档中，则 x_j 为 1，否则 x_j 为 0。注意：这种表示损失了词的所有序信息，因而称作词袋(bag)。

训练后， \hat{p}_{ij} 估计第 j 个词出现在文档类型 i 中的概率。其概率在不同类中都类似的词不传递多少信息。词要成为有用的，我们希望它的概率在一个(或少数)类上高，而在其他类上都低。我们将在第 6 章讨论这种类型的特征选择(feature selection)。文档分类应用的另一个例子是垃圾邮件过滤(spam filtering)，其中邮件被分为垃圾邮件和正常邮件两类。同样，在生物信息学中，无论碱基对还是氨基酸，输入通常也都是离散项的序列。

在一般情况下，我们假定 x_j 选自集合 $\{v_1, v_2, \dots, v_{n_j}\}$ 。我们定义新的 0/1 哑变量

$$z_{jk}^i = \begin{cases} 1 & \text{如果 } x_j^i = v_k \\ 0 & \text{否则} \end{cases}$$

令 p_{ijk} 表示 x_j 属于类 C_i 、取值为 v_k 的概率

$$p_{ijk} \equiv p(z_{jk} = 1 | C_i) = p(x_j = v_k | C_i)$$

如果属性是独立的，则我们有

$$p(\mathbf{x} | C_i) = \prod_{j=1}^d \prod_{k=1}^{n_j} p_{ijk}^{z_{jk}^i} \quad (5.32)$$

于是, 判别式函数为

$$g_i(\mathbf{x}) = \sum_j \sum_k z_{jk} \log p_{ijk} + \log P(C_i) \quad (5.33)$$

p_{ijk} 的最大似然估计为

$$\hat{p}_{ijk} = \frac{\sum_i z_{jk} r_i^t}{\sum_i r_i^t} \quad (5.34)$$

可以将它们代入式(5.33)中, 得到判别式。

5.8 多元回归

在多元线性回归(multivariate linear regression)中, 假定数值输出 r 为线性函数, 即一些输入变量 x_1, \dots, x_d 和噪声的加权和。实际上, 在统计学中, 这称为多元(multiple)回归; 当存在多个输出时, 统计学家使用术语 multivariate(多元)。多元线性模型是

$$r^t = g(\mathbf{x}^t | w_0, w_1, \dots, w_d) + \varepsilon = w_0 + w_1 x_1^t + w_2 x_2^t + \dots + w_d x_d^t + \varepsilon \quad (5.35)$$

与一元情况相同, 我们假设 ε 是正态的, 具有 0 均值和常数方差。最大化该似然等价于最小化平方误差之和:

$$E(w_0, w_1, \dots, w_d | X) = \frac{1}{2} \sum_t (r^t - w_0 - w_1 x_1^t - w_2 x_2^t - \dots - w_d x_d^t)^2 \quad (5.36)$$

对参数 $w_j (j=0, \dots, d)$ 求导, 我们得到正规方程(normal equation):

$$\sum_t r^t = N w_0 + w_1 \sum_t x_1^t + w_2 \sum_t x_2^t + \dots + w_d \sum_t x_d^t \quad (5.37)$$

$$\sum_t x_1^t r^t = w_0 \sum_t x_1^t + w_1 \sum_t (x_1^t)^2 + w_2 \sum_t x_1^t x_2^t + \dots + w_d \sum_t x_1^t x_d^t$$

$$\sum_t x_2^t r^t = w_0 \sum_t x_2^t + w_1 \sum_t x_1^t x_2^t + w_2 \sum_t (x_2^t)^2 + \dots + w_d \sum_t x_2^t x_d^t$$

\vdots

$$\sum_t x_d^t r^t = w_0 \sum_t x_d^t + w_1 \sum_t x_d^t x_1^t + w_2 \sum_t x_d^t x_2^t + \dots + w_d \sum_t (x_d^t)^2$$

我们定义如下的向量与矩阵:

$$\mathbf{X} = \begin{bmatrix} 1 & x_1^1 & x_2^1 & \dots & x_d^1 \\ 1 & x_1^2 & x_2^2 & \dots & x_d^2 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^N & x_2^N & \dots & x_d^N \end{bmatrix}, \quad \mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_d \end{bmatrix}, \quad \mathbf{r} = \begin{bmatrix} r^1 \\ r^2 \\ \vdots \\ r^N \end{bmatrix}$$

于是, 正规方程可以写为

$$\mathbf{X}^T \mathbf{X} \mathbf{w} = \mathbf{X}^T \mathbf{r} \quad (5.38)$$

并且我们可以求解参数

$$\mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{r} \quad (5.39)$$

这个方法与我们在单个输入的多项式回归中所使用的方法一样。如果我们定义变量为 $x_1 = x, x_2 = x^2, \dots, x_k = x^k$, 则两个问题是一样的。这也提示我们, 如果需要的话, 如何做多元多项式回归(multivariate polynomial regression)(参见习题 7)。但是, 除非 d 很小, 否则在多元回归中, 我们很少使用比线性更高阶的多项式。

实际上,使用输入的高阶项作为附加的输入只是一种可能的方法,我们可以使用基函数(basis function)定义原始输入的任意非线性函数。例如,如果我们相信这种变换是有用的,可以定义新的输入 $x_2 = \sin(x)$, $x_3 = \exp(x^2)$ 。然后,在这种新扩展的空间中使用线性模型对应于原空间中的非线性模型。同样的计算仍然有效,我们只需在使用基函数之后用数据矩阵替换 \mathbf{X} 。正如我们将在各种形式下(例如,多层感知器、支持向量机、高斯过程)所看到的,这些形式经常用这样的方法推广线性模型。

线性模型的一个优点是,回归之后,观察 $w_j (j = 1, \dots, d)$ 的值,我们可以提取知识:首先,观察 w_j 的符号,我们就知道 x_j 对输出结果的影响是正的还是负的;第二,如果所有的 x_j 具有相同的值域,则通过观测 w_j 的绝对值,我们可以知道特征的重要性,并按重要性为特征定秩,甚至可以去掉那些 w_j 接近于 0 的特征。

当有多个输出时,可以等价地定义一组独立的单输出回归问题。

5.9 注释

一本更新我们的线性代数知识的好书是 Strang 1988。Harville 1997 是另外一本很好的书,它从统计学的角度处理矩阵代数。

用多元数据的一个不便之处是,当维数很大时,不能够进行可视分析。统计学文献中已经提出了一些方法,来显示多元数据;Rencher 1995 给出了综述。一种可能的方法是两两变量绘制二元散点图:如果数据是多元正态的,则任意两个变量的图应该是大致上线性的;这可以用作多元正态性的视觉检测。我们将在第 6 章中讨论的另外一个可能的方法是把它们投影到一维或两维上,并在那里显示。

模式识别的大部分工作都是在假定多元正态密度上进行的。有时,这样的判别式甚至被称为贝叶斯最优分类,但是这通常是错误的;只有当密度确实是多元正态的,并且我们有足够的数据来计算正确的参数时,它才是最优的。Rencher 1995 讨论了评估多元正态性的检验,以及检查相等协方差矩阵的检验。McLachlan 1992 讨论了用多元正态分布分类,并且比较了线性和二次判别式。

多元正态分布的一个明显的约束是它不允许某些特征是离散型数据。一个具有 n 个可能值的变量可以被转化成 n 个 0/1 哑变量,但是这增加了维度。我们可以用第 6 章中介绍的方法在这个 n 维空间上进行维度归约,从而不会增加维度。对于这种混合特征参数分类,McLachlan 1992 有详细的讨论。

5.10 习题

1. 证明式(5.11)。
2. 从多元正态密度 $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ 产生一个样本,计算 \mathbf{m} 和 \mathbf{S} 并将它们与 $\boldsymbol{\mu}$ 和 $\boldsymbol{\Sigma}$ 比较。检查样本大小变化时估计的变化情况。
3. 从两个多元正态密度 $\mathcal{N}(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i) (i = 1, 2)$ 产生样本,并对表 5-1 中的四种情况计算贝叶斯最优判别式。
4. 对于两类问题,针对表 5-1 中高斯密度的四种情况,推导:

$$\log \frac{P(C_1 | \mathbf{x})}{P(C_2 | \mathbf{x})}$$

5. 使用高斯密度的另一种可能方法是令它们都是对角的，但允许它们不同。为这种情况推导判别式。
6. 假设在二维空间有两个具有相同均值的类。可以定义何种类型的边界？
7. 假设我们有两个变量 x_1 和 x_2 ，并且我们想对它们做二次拟合，即

$$f(x_1, x_2) = w_0 + w_1 x_1 + w_2 x_2 + w_3 x_1 x_2 + w_4 (x_1)^2 + w_5 (x_2)^2$$

给定样本 $X = \{x_1^t, x_2^t, r^t\}$ ，如何找到 $w_i (i=0, \dots, 5)$ ？

8. 在回归中看到拟合一个二次模型等价于用对应于输入平方的附加输入拟合一个线性模型。对于分类我们也能这样做吗？
9. 在文档聚类中，通过考虑上下文，例如考虑诸如“cocktail party”与“party elections”中的词对，可以减少二义性。讨论这一点如何实现。

106

5.11 参考文献

- Duda, R. O., P. E. Hart, and D. G. Stork. 2001. *Pattern Classification*, 2nd ed. New York: Wiley.
- Friedman, J. H. 1989. “Regularized Discriminant Analysis.” *Journal of American Statistical Association* 84: 165–175.
- Harville, D. A. 1997. *Matrix Algebra from a Statistician's Perspective*. New York: Springer.
- Manning, C. D., and H. Schütze. 1999. *Foundations of Statistical Natural Language Processing*. Cambridge, MA: MIT Press.
- McLachlan, G. J. 1992. *Discriminant Analysis and Statistical Pattern Recognition*. New York: Wiley.
- Rencher, A. C. 1995. *Methods of Multivariate Analysis*. New York: Wiley.
- Strang, G. 1988. *Linear Algebra and its Applications*, 3rd ed. New York: Harcourt Brace Jovanovich.

107
108

维 度 归 约

任何分类器和回归器的复杂度都依赖于输入的数量。这决定了时间和空间的复杂度以及训练这样的分类器和回归器所需要的训练样例数量。本章，我们讨论特征选择和特征提取方法。前者选取重要特征子集并剪掉其余特征，而后者由原始输入形成较少的新特征。

6.1 引言

在一个应用中，无论它是分类还是回归，我们确信含有信息的观测数据都被用作输入并且输入到系统中来做决策。理想情况下，我们不应该将特征选择或特征提取作为一个单独的进程；分类方法（或回归方法）应该能够利用任何必要的特征，而丢弃不相关的特征。然而，有许多原因使得我们对把降维作为一个单独的预处理步骤感兴趣：

- 在大多数学习算法中，复杂度依赖于输入的维度 d 和数据样本的规模 N ，并且为了减少存储量和计算时间，我们对降低问题的维度感兴趣。降低 d 也降低了检验时推理算法的复杂度。
- 当一个输入被认定并不必要时，我们就节省了提取它的开销。
- 较简单的模型在小数据集上更为鲁棒。较简单的模型具有较小的方差，也就是说，它们的变化更少地依赖于样本的特殊性，包括噪声、离群点等。
- 当数据能够用较少特征解释时，我们就能够更好地理解解释数据的过程，这使得我们能够提取知识。
- 当数据可以用少数维表示而不丢失信息时，我们可以对数据绘图，并可视化地分析它的结构和离群点。

降低维度的主要方法有两种：特征选择和特征提取。在特征选择（feature selection）中，我们感兴趣的是从 d 维中找出为我们提供最多信息的 k 个维，并且丢弃其他的 $(d - k)$ 个维。作为一种特征选择方法，我们将要讨论子集选择（subset selection）。

在特征提取（feature extraction）中，我们感兴趣的是找出 k 个维的新集合，这些维是原来 d 个维的组合。这些方法可以是监督的，也可以是非监督的，这取决于它们是否使用输出信息。最著名和最广泛使用的特征提取方法是主成分分析（PCA）和线性判别分析（LDA）。它们都是线性投影方法，分别是非监督和监督的。PCA 与其他两种非监督的线性投影方法有许多相似之处，我们也将讨论这两种方法——因子分析（FA）和多维定标（MDS）。作为非线性维度归约的例子，我们将考察等距特征映射（Isometric feature mapping, Isomap）和局部线性嵌入（Locally Linear Embedding, LLE）。

6.2 子集选择

在子集选择（subset selection）中，我们对发现特征集中的最佳子集感兴趣。最佳子集包

含的维最少，而它们对正确率的贡献最大。我们丢弃剩余的不重要的维。使用一个合适的误差函数，最佳子集在回归和分类问题中都可以使用。 d 个变量有 2^d 个可能子集，但是除非 d 很小，否则我们不能对所有子集进行检验。我们使用启发式的方法，在合理的（多项式）时间内得到一个合理的（但不是最优的）解。

110

有两种方法：在向前选择(forward selection)中，我们从空集开始，逐个添加它们，每次添加一个降低误差最多的变量，直到进一步的添加不会降低误差(或降低很少)。在向后选择(backward selection)中，我们从所有变量开始，逐个排除它们，每次排除一个降低误差最多(或提高很少)的变量，直到进一步的排除会显著提高误差。在这两种情况下，误差检测都应该在不同于训练集的验证集上进行，因为我们想要检验泛化准确率。使用更多的特征，我们一般会有更低的训练误差，但是不一定有更低的验证误差。

我们用 F 表示输入维的特征 $x_i (i=1, \dots, d)$ 的集合， $E(F)$ 表示当我们只使用 F 中的输入时，在验证样本上出现的误差。依赖于应用，误差或者是均方差误差，或者是误分类错误。

在顺序向前选择(sequential forward selection)中，我们从 $F = \emptyset$ 开始。每一步中，我们针对所有可能的 x_i ，训练我们的模型并在验证集上计算 $E(F \cup x_i)$ 。然后，我们选择导致最小误差的输入 x_j

$$j = \arg \min_i E(F \cup x_i) \quad (6.1)$$

并且我们

$$\text{将 } x_j \text{ 添加到 } F \text{ 中, 如果 } E(F \cup x_j) < E(F) \quad (6.2)$$

如果添加任何特征都不会减少 E ，则我们停止。如果误差降低太小，我们甚至可以决定提前停止；这里存在一个用户定义的阈值，依赖于应用约束以及错误和复杂度的折中。增加另外一个特征带来观测该特征的开销，也会使分类器/回归器更加复杂。

这样的过程也许开销很大，因为将 d 维减少到 k 维，我们需要训练和测试系统 $d + (d-1) + (d-2) + \dots + (d-k)$ 次，其复杂度为 $\mathcal{O}(d^2)$ 。这是一个局部搜索过程，并且不能保证找到最佳子集，即导致最小误差的最小子集。例如， x_i 和 x_j 本身可能不好，但是合起来却可能会把误差降低很多。但是该算法很贪婪，逐个增加特征，因此它也许不能发现 x_i 与 x_j 的并。以更多计算为代价，一次增加 m 个而不是一个特征是可能的。我们还可以在当前添加之后回溯并且检查以前添加的哪个特征可以去掉，这增大了搜索空间但是也增加了复杂度。在浮动搜索(floating search)方法中(Pudil、Novovičová 和 Kittler 1994)，每一步还可以改变增加和去掉的特征数量。

111

在顺序向后选择(sequential backward selection)中，我们从包括所有特征的 F 开始，并且执行类似的过程，但是与添加相反，我们从 F 中去掉一个特征，并且是去掉导致误差最小的那个

$$j = \arg \min_i E(F - x_i) \quad (6.3)$$

而我们

$$\text{从 } F \text{ 中去掉 } x_j, \text{ 如果 } E(F - x_j) < E(F) \quad (6.4)$$

如果去掉特征不能降低误差我们就停止。为了降低复杂度，我们可能也会决定去掉一个特征，如果它的去掉只引起很轻微的误差增加。

向前搜索的所有可能变体对于向后搜索也是可行的。向后搜索与向前搜索具有相同的复杂度。但是，训练具有较多特征的系统比训练具有较少特征的系统开销更大，并且如果我们

预料有许多无用特征时，向前搜索更可取。

子集选择是监督的，因为输出被回归器或分类器用作计算误差，但是它可以用于任何回归和分类方法。在多元正态分类的特殊情况下，如果原来的 d 维类密度是多元正态的，则其任意子集也是多元正态的，并且仍然可以使用参数分类，并具有用 $k \times k$ 维协方差矩阵代替 $d \times d$ 维协方差矩阵的优点。

在像人脸识别这样的应用中，特征选择不是降维的好方法，因为个体像素本身并不携带很多识别信息；携带脸部识别信息的是许多像素值的组合。这可以通过我们后边将要讨论的特征提取方法来做。

112

6.3 主成分分析

在投影方法中，我们感兴趣的是找到一个从原 d 维输入空间到新的 ($k < d$) 维空间的、具有最小信息损失的映射。 \mathbf{x} 在方向 \mathbf{w} 上的投影为

$$\mathbf{z} = \mathbf{w}^T \mathbf{x} \quad (6.5)$$

主成分分析 (principal components analysis, PCA) 是一种非监督方法，因为它不使用输出信息；需要最大化的是方差。主成分是这样的 \mathbf{w}_1 ，样本投影到 \mathbf{w}_1 上之后被广泛散布，使得样本点之间的差别变得最明显。为了得到唯一解并且使该方向成为最重要因素，我们要求 $\|\mathbf{w}_1\| = 1$ 。从式(5.14)我们知道，如果 $\mathbf{z}_1 = \mathbf{w}_1^T \mathbf{x}$ ，并且 $\text{Cov}(\mathbf{x}) = \Sigma$ ，则

$$\text{Var}(\mathbf{z}_1) = \mathbf{w}_1^T \Sigma \mathbf{w}_1$$

我们寻找 \mathbf{w}_1 ，使得 $\text{Var}(\mathbf{z}_1)$ 受限于约束 $\mathbf{w}_1^T \mathbf{w}_1 = 1$ 最大化。将这写成拉格朗日问题，我们得到

$$\max_{\mathbf{w}_1} \mathbf{w}_1^T \Sigma \mathbf{w}_1 - \alpha (\mathbf{w}_1^T \mathbf{w}_1 - 1) \quad (6.6)$$

关于 \mathbf{w}_1 求导并令它等于 0，我们有

$$2\Sigma \mathbf{w}_1 - 2\alpha \mathbf{w}_1 = 0, \quad \text{因此 } \Sigma \mathbf{w}_1 = \alpha \mathbf{w}_1$$

如果 \mathbf{w}_1 是 Σ 的特征向量， α 是对应的特征值，则上式成立。因为我们有

$$\mathbf{w}_1^T \Sigma \mathbf{w}_1 = \alpha \mathbf{w}_1^T \mathbf{w}_1 = \alpha$$

为了使方差最大，我们选择具有最大特征值的特征向量。因此，主成分是输入样本的协方差矩阵的具有最大特征值 $\lambda_1 = \alpha$ 的本征向量。

第二个主成分 \mathbf{w}_2 也应该最大化方差，具有单位长度，并且与 \mathbf{w}_1 正交。后者的要求是使得投影后 $\mathbf{z}_2 = \mathbf{w}_2^T \mathbf{x}$ 与 \mathbf{z}_1 不相关。对于第二个主成分，我们有

$$\max_{\mathbf{w}_2} \mathbf{w}_2^T \Sigma \mathbf{w}_2 - \alpha (\mathbf{w}_2^T \mathbf{w}_2 - 1) - \beta (\mathbf{w}_2^T \mathbf{w}_1 - 0) \quad (6.7) \quad 113$$

关于 \mathbf{w}_2 求导并令它等于 0，我们有

$$2\Sigma \mathbf{w}_2 - 2\alpha \mathbf{w}_2 - \beta \mathbf{w}_1 = 0 \quad (6.8)$$

用 \mathbf{w}_1^T 左乘，我们得到

$$2\mathbf{w}_1^T \Sigma \mathbf{w}_2 - 2\alpha \mathbf{w}_1^T \mathbf{w}_2 - \beta \mathbf{w}_1^T \mathbf{w}_1 = 0$$

注意 $\mathbf{w}_1^T \mathbf{w}_2 = 0$ 。 $\mathbf{w}_1^T \Sigma \mathbf{w}_2$ 是标量，等于它的转置 $\mathbf{w}_2^T \Sigma \mathbf{w}_1$ ，这里 \mathbf{w}_1 是 Σ 的主特征向量， $\Sigma \mathbf{w}_1 = \lambda_1 \mathbf{w}_1$ ，因此

$$\mathbf{w}_1^T \Sigma \mathbf{w}_2 = \mathbf{w}_2^T \Sigma \mathbf{w}_1 = \lambda_1 \mathbf{w}_2^T \mathbf{w}_1 = 0$$

于是 $\beta = 0$ ，并且式(6.8)可以简化为

$$\Sigma \mathbf{w}_2 = \alpha \mathbf{w}_2$$

这表明 \mathbf{w}_2 应该是 Σ 的特征向量, 具有第二大特征值 $\lambda_2 = \alpha$ 。类似地, 我们可以证明其他维被具有递减的本征值的本征向量给出。

因为 Σ 是对称的, 因此对于两个不同的特征值, 特征向量是正交的。如果 Σ 是正定的 (对于所有的非空 \mathbf{x} , $\mathbf{x}^T \Sigma \mathbf{x} > 0$), 则它的所有特征值都是正的。如果 Σ 是奇异的, 则它的秩 (有效维数) 为 k , 并且 $k < d$, $\lambda_i (i = k+1, \dots, d)$ 均为 0 (λ_i 以递减序排序)。 k 个具有非零本征值的特征向量是约化空间的维。第一个特征向量 (具有最大特征值的向量) \mathbf{w}_1 (即为主成分) 贡献了方差的最大部分, 第二个贡献了方差的第二大部分, 依此类推。

我们定义

$$\mathbf{z} = \mathbf{W}^T (\mathbf{x} - \mathbf{m}) \quad (6.9)$$

其中 \mathbf{W} 的 k 列是 \mathbf{S} 的 k 个主特征向量, 也是 Σ 的估计。我们在 \mathbf{x} 投影前减去均值 \mathbf{m} , 将数据在原点中心化。该线性变换后, 我们得到 k 维空间, 它的维是特征向量, 并且在这些新维上的方差等于特征值 (见图 6-1)。为了规范化方差, 我们可以除以特征值的平方根。

114

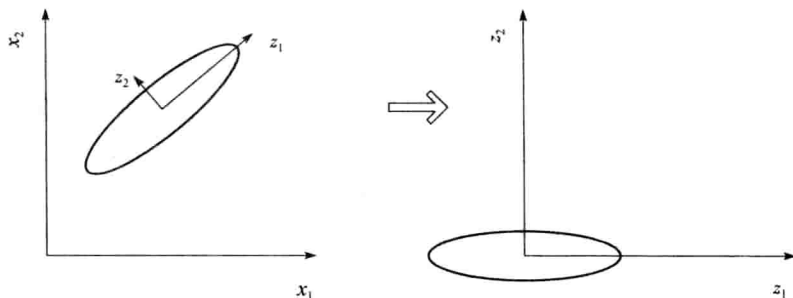


图 6-1 主成分分析使样本中心化, 然后旋转坐标轴与最大方差方向一致。如果 z_2 上的方差太小, 则可以忽略它, 并且我们得到从二维到一维的维度归约

让我们来看另一种推导: 我们想要找到一个矩阵 \mathbf{W} , 使得当我们有 $\mathbf{z} = \mathbf{W}^T \mathbf{x}$ (假设不失一般性, \mathbf{x} 已经被中心化) 时, 我们将得到 $\text{Cov}(\mathbf{z}) = \mathbf{D}'$, 其中 \mathbf{D}' 是任意对角阵; 也就是说, 我们希望得到不相关的 z_i 。

如果我们建立一个 $d \times d$ 的矩阵 \mathbf{C} , 其第 i 列是 \mathbf{S} 的规范化的特征向量 \mathbf{c}_i , 则 $\mathbf{C}^T \mathbf{C} = \mathbf{I}$, 并且

$$\begin{aligned} \mathbf{S} &= \mathbf{S} \mathbf{C} \mathbf{C}^T \\ &= \mathbf{S} (\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_d) \mathbf{C}^T \\ &= (\mathbf{S} \mathbf{c}_1, \mathbf{S} \mathbf{c}_2, \dots, \mathbf{S} \mathbf{c}_d) \mathbf{C}^T \\ &= (\lambda_1 \mathbf{c}_1, \lambda_2 \mathbf{c}_2, \dots, \lambda_d \mathbf{c}_d) \mathbf{C}^T \\ &= \lambda_1 \mathbf{c}_1 \mathbf{c}_1^T + \dots + \lambda_d \mathbf{c}_d \mathbf{c}_d^T \\ &= \mathbf{C} \mathbf{D} \mathbf{C}^T \end{aligned} \quad (6.10)$$

其中 \mathbf{D} 是对角矩阵, 其对角线元素是特征值 $\lambda_1, \dots, \lambda_d$ 。这称为 \mathbf{S} 的谱分解 (spectral decomposition)。由于 \mathbf{C} 是正交的, 并且 $\mathbf{C} \mathbf{C}^T = \mathbf{C}^T \mathbf{C} = \mathbf{I}$, 我们可以对上式左乘以 \mathbf{C}^T , 右乘以 \mathbf{C} , 得到

$$\mathbf{C}^T \mathbf{S} \mathbf{C} = \mathbf{D} \quad (6.11)$$

我们知道如果 $\mathbf{z} = \mathbf{W}^T \mathbf{x}$, 则 $\text{Cov}(\mathbf{z}) = \mathbf{W}^T \mathbf{S} \mathbf{W}$, 我们希望它等于一个对角矩阵。于是, 从式(6.11)我们看到, 可以令 $\mathbf{W} = \mathbf{C}$ 。

115

让我们看一个例子, 以便得到一些直观体验(Rencher 1995): 假设我们有一个班学生的五门课程的成绩, 我们希望对这些学生排序。也就是说, 我们希望把这些数据投影到一个维上, 使得这些数据点之间的差别最明显。我们可以用 PCA。具有最大特征值的特征向量是最大方差的方向, 也就是学生最为分散的方向。这样做比计算平均值好, 因为我们考虑了方差的相关性和区别。

实践中, 即使所有特征值都大于 0, 但是如果 $|\mathbf{S}|$ 很小 (注意 $|\mathbf{S}| = \prod_{i=1}^d \lambda_i$), 那么我们知道, 某些特征值对方差影响很小, 并且可以丢弃。因此, 我们考虑例如贡献 90% 以上方差的前 k 个主要成分。当 λ_i 降序排列时, 由前 k 个主要成分贡献的方差比例 (proportion of variance) 为

$$\frac{\lambda_1 + \lambda_2 + \cdots + \lambda_k}{\lambda_1 + \lambda_2 + \cdots + \lambda_k + \cdots + \lambda_d}$$

如果维是高度相关的, 则只有很少一部分特征向量具有较大的本征值, k 远比 d 小, 并且可能得到很大的维度归约。在许多图像和语音处理任务中, 通常是这种情况, 其中(时间或空间)邻近的输入是高度相关的。如果维之间互不相关, k 将与 d 一样大, 通过 PCA 就没有收益。

斜坡图 (scree graph) 是把贡献的方差作为特征向量编号的函数的图形 (见图 6-2)。通过目视分析, 我们也可以确定 k 。在“拐点”处, 增加其他特征向量不会显著地增加贡献的方差。

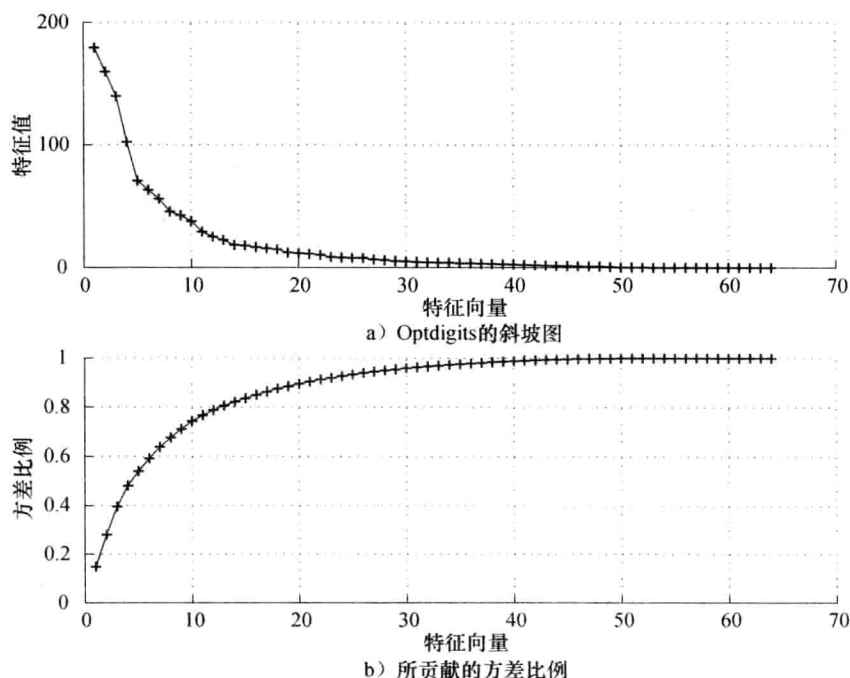


图 6-2 a) 斜坡图。b) 对于取自 UCI 知识库的 Optdigits 数据集, 显示所贡献的方差比例。Optdigits 是手写数字数据集, 具有 10 个类和 64 维输入。前 20 个特征向量贡献了 90% 的方差

另一个可能的的方法是忽略那些特征值小于平均输入方差的特征向量。给定 $\sum_i \lambda_i = \sum_i s_i^2$ (等于矩阵 S 的迹, 记作 $\text{tr}(S)$), 平均特征值等于平均输入方差。当我们仅保留特征值大于平均特征值的特征向量时, 我们仅保留了那些其方差大于平均输入方差的特征向量。

如果原 x_i 维的方差变化显著, 则它们对主成分方向的影响比相关性大。因此, 一个公共过程是在使用 PCA 之前对数据进行预处理, 使得每个维都具有 0 均值和单位方差。或者, 为了使协方差而不是个体方差起作用, 我们可以使用协相关矩阵 R 而不是协方差矩阵 S 的本征向量。

PCA 解释方差并对离群点很敏感: 少量远离中心的点对方差有很大影响, 从而也对特征向量有很大影响。鲁棒的估计 (Robust estimation) 方法允许计算离群点存在时的参数。一种简单的方法是计算数据点的马氏距离, 丢弃那些远离的孤立数据点。

如果前两个主成分贡献方差的很大百分比, 则我们可以做目视分析: 我们可以在这个二维空间绘制数据 (见图 6-3), 目视地搜索结构、组、离群点、正态性等。相对于原来的任何两个变量的图, 该图对样本给出了更好的图形描述。通过观察主成分的维, 我们还可以试着揭示一些有意义的描述数据的潜在变量。例如, 在图像应用方面, 输入是图像, 特征向量可以显示为图像, 并且可以看作重要特征的模板; 它们常常被形象地称为“特征面孔” (eigen-face)、“特征数字” (eigendigit) 等 (Turk 和 Pentland 1991)。

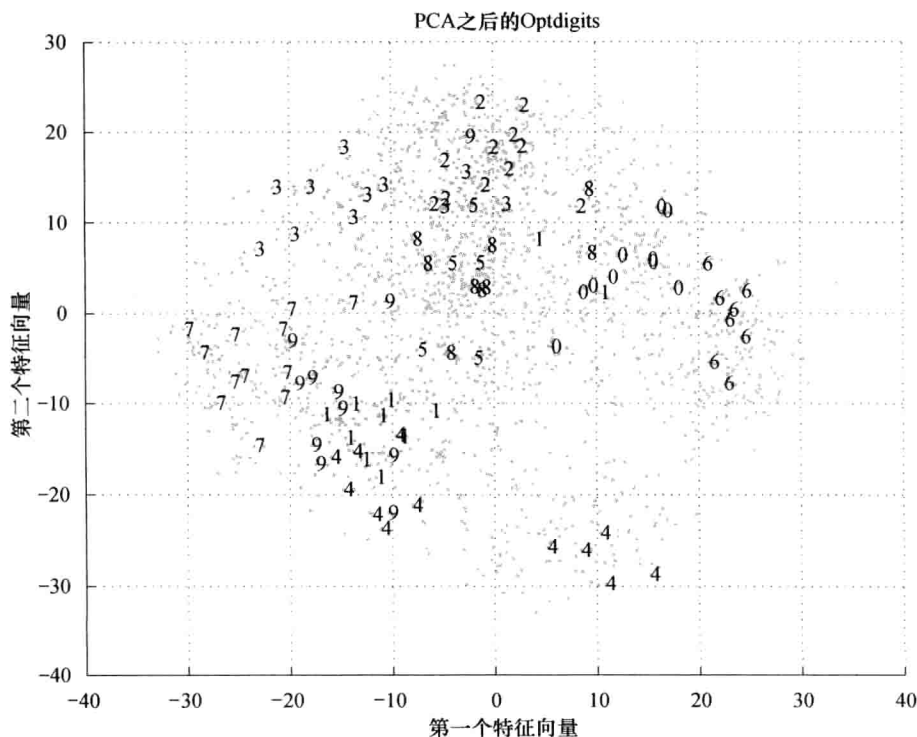


图 6-3 绘制在两个主成分空间的 Optdigits 数据。只显示了 100 个数据点的标号, 以便最小化墨噪比 (ink-to-noise ratio)

当 d 很大时, 计算、存储、处理 S 都很烦琐。我们可以直接从数据计算特征向量、特征值, 而不必显式地计算协方差矩阵(Chatfield 和 Collins 1980)。

118

从式(5.15)我们知道, 如果 $\mathbf{x} \sim \mathcal{N}_d(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, 则投影后 $\mathbf{W}^T \mathbf{x} \sim \mathcal{N}_k(\mathbf{W}^T \boldsymbol{\mu}, \mathbf{W}^T \boldsymbol{\Sigma} \mathbf{W})$ 。如果样本是 d 元正态的, 则它投影到 k 元正态上, 允许我们在很有希望的、低得多的维空间进行参数判别分析。因为 z_j 是不相关的, 因此新的协方差矩阵将是对角的。如果它们被规范化, 具有单位方差, 则可以在这个新空间使用欧氏距离, 导出简单的分类器。

实例 \mathbf{x}' 投影到 \mathbf{z} -空间

$$\mathbf{z}' = \mathbf{W}^T(\mathbf{x}' - \boldsymbol{\mu})$$

当 \mathbf{W} 是正交矩阵使得 $\mathbf{W}\mathbf{W}^T = \mathbf{I}$ 时, 它可以逆投影到原来的空间

$$\hat{\mathbf{x}}' = \mathbf{W}\mathbf{z}' + \boldsymbol{\mu}$$

$\hat{\mathbf{x}}'$ 是 \mathbf{x}' 从它在 \mathbf{z} -空间中的表示的重构。我们知道, 在所有正交线性投影中, PCA 最小化重构误差(reconstruction error)。重构误差是实例与它的从低维空间重构之间的距离:

$$\sum_i \|\hat{\mathbf{x}}' - \mathbf{x}'\|^2 \quad (6.12)$$

重构误差取决于考虑了多少个主成分。在视觉识别应用中, 例如, 人脸识别, 显示 $\hat{\mathbf{x}}'$ 使得我们能够可视化地检查 PCA 期间的信息损失。

PCA 是非监督的, 并且不利用输出信息。它是一个一组(one-group)过程。然而, 在分类情况下会有很多组, Karhunen-Loève 扩展(Karhunen-Loève expansion)允许利用类信息; 例如, 我们不是使用整个样本的协方差矩阵, 而是估计类的协方差矩阵, 取它们的平均(用先验加权)作为协方差矩阵, 并使用它的特征向量。

在公共主成分(common principal component)中(Flury 1988), 我们假设对于每个类, 主成分都是相同的, 但是对于不同的类, 这些成分的分量不同:

$$\mathbf{S}_i = \mathbf{C}\mathbf{D}_i\mathbf{C}^T$$

这种方法允许汇聚数据, 并且是一种正则化方法, 它的复杂度比所有类的公共协方差矩阵的方法小, 同时仍允许 \mathbf{S}_i 存在差别。一种相关的方法是柔性判别分析(flexible discriminant analysis)(Hastie、Tibshirani 和 Buja 1994), 它将数据线性投影到所有特征都是不相关的低维空间, 再用最小距离分类。

119

6.4 因子分析

在 PCA 中, 从原来的维 $x_i (i=1, \dots, d)$, 我们形成一个新的变量集 \mathbf{z} , 它是 x_i 的线性组合:

$$\mathbf{z} = \mathbf{W}^T(\mathbf{x} - \boldsymbol{\mu})$$

在因子分析(factor analysis, FA)中, 我们假定有一个不可观测的潜在因子(latent factor) $z_j (j=1, \dots, k)$ 的集合, 它在组合时生成 \mathbf{x} 。因此, 与 PCA 的方向相反(见图 6-4), 其目标是通过较少数量的因子刻画观测变量之间的依赖性。

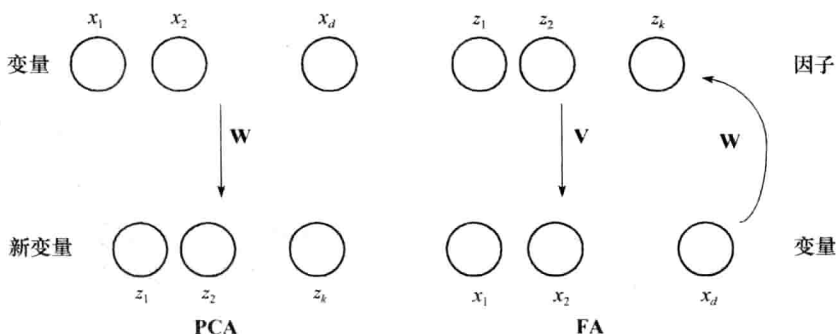


图 6-4 主成分分析方法产生新的变量，它们是原输入变量的线性组合。然而，在因子分析，我们假定存在一些因子，它们线性组合时产生输入变量

假设有一个变量组，它们之间具有高度相关性，而与其他所有变量具有很低的相关性。那么可能存在于一个简单的潜在因子给出这些变量的起源。如果其他变量能够类似地分成子集，则少数因子就能够代表这些变量组。虽然因子分析总是把变量划分成因子簇，但是因子是否意味着什么，或是否真的存在，仍然是一个悬而未决的问题。

像 PCA 一样，FA 也是一个一组过程，并且是非监督的。目标是在一个更小的维空间中对数据建模而不丢失信息。在 FA 中，这用变量之间的相关性度量。

正如在 PCA 中一样，我们有样本 $\mathcal{X} = \{\mathbf{x}^i\}_i$ ，取自某个未知的概率密度，其中 $E[\mathbf{x}] = \boldsymbol{\mu}$ ， $\text{Cov}(\mathbf{x}) = \boldsymbol{\Sigma}$ 。我们假定因子是单位正态的 $E[z_j] = 0$ ， $\text{Var}(z_j) = 1$ ，并且是不相关的 $\text{Cov}(z_i, z_j) = 0$ ， $i \neq j$ 。为了说明什么是不能由因子解释的，每个输入存在一个附加的源，记作 ε_i 。假定它具有 0 均值 $E[\varepsilon_i] = 0$ ，和某个未知的方差 $\text{Var}(\varepsilon_i) = \Psi_i$ 。这些特殊的源之间是不相关的 $\text{Cov}(\varepsilon_i, \varepsilon_j) = 0$ ， $i \neq j$ ，并且与因子也是不相关的 $\text{Cov}(\varepsilon_i, z_j) = 0$ ， $\forall i, j$ 。

FA 假定每个输入维 $x_i (i=1, \dots, d)$ 可以写成 $k < d$ 个因子 $z_j (j=1, \dots, k)$ 的加权和，加上残差项(见图 6-5)：

$$x_i - \mu_i = v_{i1}z_1 + v_{i2}z_2 + \dots + v_{ik}z_k + \varepsilon_i, \quad \forall i = 1, \dots, d$$

$$x_i - \mu_i = \sum_{j=1}^k v_{ij}z_j + \varepsilon_i \quad (6.13)$$

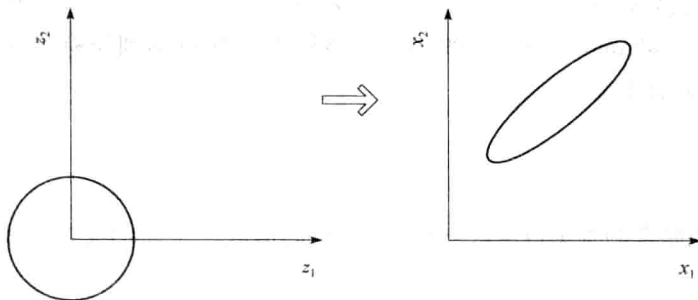


图 6-5 因子是独立的、单位正态的，它们被延伸、旋转和变换，以成为输入

这可以写成向量矩阵形式

$$\mathbf{x} - \boldsymbol{\mu} = \mathbf{V}\mathbf{z} + \boldsymbol{\varepsilon} \quad (6.14)$$

其中 V 是 $d \times k$ 权重矩阵, 称作因子载荷 (factor loading)。从现在开始, 不失一般性, 我们将假设 $\mu = 0$; 我们总能在投影后加上 μ 。给定 $\text{Var}(z_j) = 1$ 和 $\text{Var}(\varepsilon_i) = \Psi_i$

$$\text{Var}(x_i) = v_{i1}^2 + v_{i2}^2 + \cdots + v_{ik}^2 + \Psi_i \quad (6.15)$$

$\sum_{j=1}^k v_{ij}^2$ 是公共因子贡献的部分, 而 Ψ_i 是针对 x_i 的方差。

用向量矩阵形式, 我们有

$$\Sigma = \text{Cov}(\mathbf{x}) = \text{Cov}(\mathbf{V}\mathbf{z} + \boldsymbol{\varepsilon}) \quad (6.16)$$

$$= \text{Cov}(\mathbf{V}\mathbf{z}) + \text{Cov}(\boldsymbol{\varepsilon})$$

$$= \mathbf{V}\text{Cov}(\mathbf{z})\mathbf{V}^T + \boldsymbol{\Psi}$$

$$= \mathbf{V}\mathbf{V}^T + \boldsymbol{\Psi} \quad (6.17)$$

120
1
121

其中 $\boldsymbol{\Psi}$ 是对角矩阵, Ψ_i 在对角线上。由于因子是不相关的、单位正态的, 因此我们有 $\text{Cov}(\mathbf{z}) = \mathbf{I}$ 。例如, 对于两个因子

$$\text{Cov}(x_1, x_2) = v_{11}v_{21} + v_{12}v_{22}$$

如果 x_1 和 x_2 的协方差高, 则它们通过一个因子相关。如果它是第一个因子, 则 v_{11} 和 v_{21} 都高; 如果它是第二个因子, 则 v_{12} 和 v_{22} 都高。在这两种情况下, 和 $v_{11}v_{21} + v_{12}v_{22}$ 都将会高。如果该协方差低, 则 x_1 和 x_2 依赖于不同的因子, 并且在和式的乘积中, 一项高而另一项低, 而它们的和低。

我们看到

$$\text{Cov}(x_1, z_2) = \text{Cov}(v_{12}z_2, z_2) = v_{12}\text{Var}(z_2) = v_{12}$$

因此, $\text{Cov}(\mathbf{x}, \mathbf{z}) = \mathbf{V}$, 并且我们看到载荷用因子表示变量之间的相关性。

给定 Σ 的估计 S , 我们希望求解 V 和 $\boldsymbol{\Psi}$, 满足

$$S = \mathbf{V}\mathbf{V}^T + \boldsymbol{\Psi}$$

如果只有少量因子, 即如果 V 只有少数几列, 则因为 V 是 $d \times k$ 而 $\boldsymbol{\Psi}$ 有 d 个值, 我们就可能有一个关于 S 的简化结构, 这样参数的数量从 d^2 减少到 $d \cdot k + d$ 。

因为 $\boldsymbol{\Psi}$ 是对角的, 因此协方差由 V 表示。注意, PCA 不允许单独的 $\boldsymbol{\Psi}$ 并且试图考虑协方差和方差。当所有的 Ψ_i 相等时, 即当 $\boldsymbol{\Psi} = \Psi\mathbf{I}$ 时, 我们得到概率 PCA (probabilistic PCA) (Tipping 和 Bishop, 1997), 而当 Ψ_i 为 0 时, 得到传统的 PCA。

现在, 让我们来看怎样找到因子载荷和特定的方差: 我们首先忽略 $\boldsymbol{\Psi}$ 。然后, 从它的谱分解, 我们知道有

$$S = \mathbf{C}\mathbf{D}\mathbf{C}^T = \mathbf{C}\mathbf{D}^{1/2}\mathbf{D}^{1/2}\mathbf{C}^T = (\mathbf{C}\mathbf{D}^{1/2})(\mathbf{C}\mathbf{D}^{1/2})^T$$

其中, 通过观察贡献的方差比例, 我们只取 k 个本征向量, 使得 C 是 $d \times k$ 的本征向量矩阵, 而 $\mathbf{D}^{1/2}$ 是 $k \times k$ 对角矩阵, 其对角线上的值是本征值的平方根。这样, 我们有

$$\mathbf{V} = \mathbf{C}\mathbf{D}^{1/2} \quad (6.18)$$

从式(6.15)我们可以得到 Ψ_i

$$\Psi_i = s_i^2 - \sum_{j=1}^k v_{ij}^2 \quad (6.19)$$

注意当 V 与任一正交矩阵 (即具有 $\mathbf{T}\mathbf{T}^T = \mathbf{I}$ 性质的矩阵) 相乘, 它就是另一个有效解, 因此解不是唯一的。

$$S = (VT)(VT)^T = VTT^TV^T = VIV^T = VV^T$$

如果 T 是正交矩阵, 则到原点的距离不变。如果 $z = Tx$, 则

$$z^T z = (Tx)^T (Tx) = x^T T^T T x = x^T x$$

乘以一个对角矩阵具有旋转坐标轴的效果, 这允许我们选择最可解释的坐标集合 (Rencher, 1995)。在两维中,

$$T = \begin{bmatrix} \cos\phi & -\sin\phi \\ \sin\phi & \cos\phi \end{bmatrix}$$

将坐标旋转 ϕ 度。有两种类型的旋转: 在正交旋转中, 旋转后因子仍然正交; 在斜旋转中, 允许因子变成相关的。旋转因子为每个变量在尽可能少的因子上给出最大载荷, 使得因子可解释。然而, 可解释性是主观的, 不应该被用来强加个人对数据的偏见。

因子分析有两种用法: 当我们找到载荷并且试图使用较少因子表示变量时, 它可以用来提取知识。当 $k < d$ 时, 它也可以被用来降低维度。我们已经看到了第一个是怎么做的。现在, 让我们看一下因子分析怎样被用来降低维度。

当我们对维度归约感兴趣时, 我们需要能够从 x_i 发现因子得分 z_j 。我们希望找到载荷 w_{ji} , 使得

$$z_j = \sum_{i=1}^d w_{ji} x_i + \varepsilon_j, \quad j = 1, \dots, k \quad (6.20)$$

其中 x_i 被中心化, 具有均值 0。在向量形式, 对于观测 t , 这可以写作

$$z^t = W^T x^t + \varepsilon, \quad \forall t = 1, \dots, N$$

这是一个线性模型, 有 d 个输入和 k 个输出。其转置可以写作

$$(z^t)^T = (x^t)^T W + \varepsilon^T, \quad \forall t = 1, \dots, N$$

给定我们有一个 N 个观测的样本, 我们记

$$Z = XW + \Xi \quad (6.21)$$

其中 Z 是 $N \times k$ 个因子, X 是 $N \times d$ 个 (中心化的) 观测, 而 Ξ 是 $N \times k$ 个 0 均值噪声。这是一个多输出的多元线性回归, 并且我们从 5.8 节知道, 可以求解 W 得到

$$W = (X^T X)^{-1} X^T Z$$

但是我们不知道 Z ; 这是我们要计算的。我们在两边同时乘以和除以 $N-1$, 得到

$$\begin{aligned} W &= (N-1)(X^T X)^{-1} \frac{X^T Z}{N-1} \\ &= \left(\frac{X^T X}{N-1} \right)^{-1} \frac{X^T Z}{N-1} \\ &= S^{-1} V \end{aligned} \quad (6.22)$$

并且把式 (6.22) 代入式 (6.21) 中, 我们记作

$$Z = XW = XS^{-1}V \quad (6.23)$$

假定 S 是非奇异的。当 x_i 被规范化具有单位方差时, 我们可以用 R 代替 S 。

对于维度归约, 除了因子的可解释性, 允许识别公共原因、简单解释、知识提取外, FA 与 PCA 相比并无优势。例如, 在语音识别, x 对应声音信号, 但是我们知道这是少数的发音器官, 即颚、舌、软腭、嘴唇和口腔 (非线性) 相互作用的结果, 它们被适当的定位用

来形成从肺部出来的气流，进而产生语音。如果语音信号可以转换到这个发音分析的空间，则语音识别就会非常容易。使用这种生成模型是当前语音识别的研究方向之一，第 16 章我们将讨论如何用图形模型来表示这种模型。

6.5 多维定标

假设我们有 N 个点，并且给定每对点之间的距离 d_{ij} , $i, j = 1, \dots, N$ 。我们不知道这些点的确切坐标，也不知道它们的维度，以及距离是怎样计算的。多维定标 (multidimensional scaling, MDS) 是把这些点映射到低维 (如二维) 空间的方法，使得它们在两维空间中的欧氏距离尽可能接近在原空间中给定的距离 d_{ij} 。这样，它需要一个从某个未知维度空间到低维空间，例如两维空间上的投影。

在典型的多维定标例子中，我们取城市之间的道路旅行距离，在应用 MDS 后，我们得到一张近似地图。这个地图被扭曲，在存在诸如高山和湖泊等地理障碍物的部分，道路旅行距离大大的偏离了直接的飞行距离 (欧氏距离)，这个地图被拉伸，以便适应更长的距离 (见图 6-6)。该地图以原点为中心，但是解仍然不是唯一的。我们可以得到任意的旋转和镜像版本。

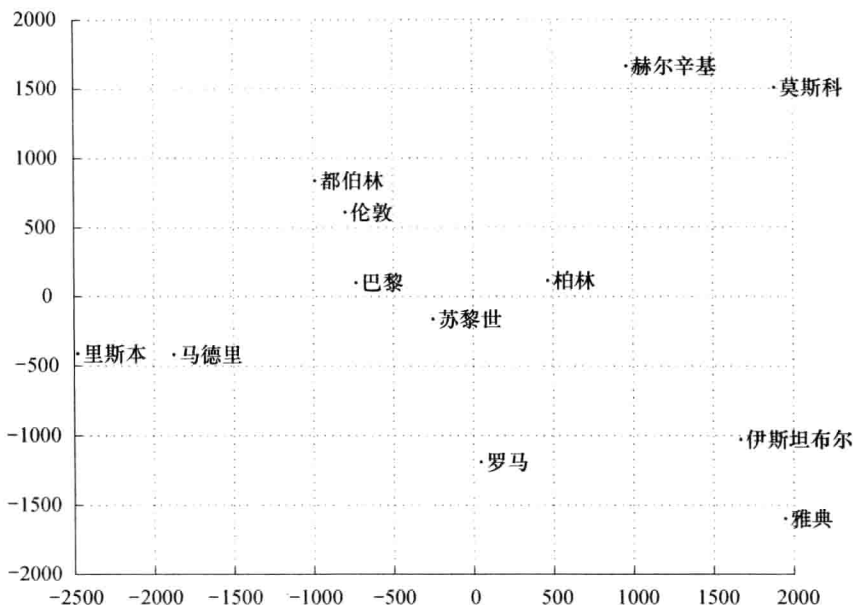


图 6-6 MDS 绘制的欧洲图。城市包括雅典、柏林、都柏林、赫尔辛基、伊斯坦布尔、里斯本、伦敦、马德里、莫斯科、巴黎、罗马和苏黎世。这些城市之间逐对道路旅行距离作为输入给出，并且 MDS 把它们放到两维的空间中，使得这些距离尽可能地被保持

可以使用 MDS 进行维度归约：通过计算 d 维 \mathbf{x} 空间的逐对欧氏距离并把它作为 MDS 的输入。然后，MDS 把它投影到较低维的空间，以保持这些距离。

假设我们有通常的样本 $\mathcal{X} = \{\mathbf{x}^i\}_{i=1}^N$ ，其中 $\mathbf{x}^i \in \mathbb{R}^d$ 。对于两个点 r 和 s ，它们之间的平方欧氏距离为

$$\begin{aligned}
 d_{rs}^2 &= \| \mathbf{x}^r - \mathbf{x}^s \|^2 = \sum_{j=1}^d (x_j^r - x_j^s)^2 = \sum_{j=1}^d (x_j^r)^2 - 2 \sum_{j=1}^d x_j^r x_j^s + \sum_{j=1}^d (x_j^s)^2 \\
 &= b_{rr} + b_{ss} - 2b_{rs}
 \end{aligned} \quad (6.24)$$

其中 b_{rs} 定义为

$$b_{rs} = \sum_{j=1}^d x_j^r x_j^s \quad (6.25)$$

为了约束这个解，我们把数据在原点中心化并假定

$$\sum_{i=1}^N x_j^i = 0, \quad \forall j = 1, \dots, d$$

然后，在 r, s 和 r, s 二者上把式(6.24)加起来，并且定义

$$T = \sum_{i=1}^N b_{ii} = \sum_i \sum_j (x_j^i)^2$$

我们得到

$$\begin{aligned}
 \sum_r d_{rs}^2 &= T + Nb_{ss} \\
 \sum_s d_{rs}^2 &= Nb_{rr} + T \\
 \sum_r \sum_s d_{rs}^2 &= 2NT
 \end{aligned}$$

当我们定义

$$d_{\cdot s}^2 = \frac{1}{N} \sum_r d_{rs}^2, \quad d_{r \cdot}^2 = \frac{1}{N} \sum_s d_{rs}^2, \quad d_{\cdot \cdot}^2 = \frac{1}{N^2} \sum_r \sum_s d_{rs}^2$$

并使用式(6.24)时，我们得到

$$b_{rs} = \frac{1}{2}(d_{r \cdot}^2 + d_{\cdot s}^2 - d_{\cdot \cdot}^2 - d_{rs}^2) \quad (6.26)$$

现在，已经计算了 b_{rs} 并已知 $\mathbf{B} = \mathbf{X}\mathbf{X}^T$ (如式 6.25 中的定义)，我们寻找一个近似。从谱分解我们知道 $\mathbf{X} = \mathbf{C}\mathbf{D}^{1/2}$ 可以被用作 \mathbf{X} 的一个近似，其中 \mathbf{C} 是矩阵，其列是 \mathbf{B} 的特征向量，而 $\mathbf{D}^{1/2}$ 是对角矩阵，其对角线是本征值的平方根。观察 \mathbf{B} 的本征值，像我们在 PCA 和 FA 中所做的那样，我们确定比 d (和 N) 低的维度 k 。我们假设 \mathbf{c}_j 是本征向量，其对应的本征值为 λ_j 。注意 \mathbf{c}_j 是 N 维的。于是，我们得到新的维

$$z_j^t = \sqrt{\lambda_j} c_j^t, \quad j = 1, \dots, k, \quad t = 1, \dots, N \quad (6.27)$$

也就是说，在标准化后，实例 t 的新坐标由本征向量 $\mathbf{c}_j (j=1, \dots, k)$ 的第 t 个元素给出。

已经证明 (Chatfield 和 Collins 1980) $\mathbf{X}\mathbf{X}^T (N \times N)$ 与 $\mathbf{X}^T \mathbf{X} (d \times d)$ 的本征值相同，并且本征向量通过一个简单的线性变换相关。这说明 PCA 做了与 MDS 相同的工作，并且代价更低。在相关矩阵而不是在协方差矩阵上做 PCA 等价于用标准化的欧氏距离做 MDS，其中每个变量有单位方差。

在一般情况下，我们希望寻找一个映射 $\mathbf{z} = \mathbf{g}(\mathbf{x} | \theta)$ ，其中 $\mathbf{z} \in \mathbb{R}^k$ ， $\mathbf{x} \in \mathbb{R}^d$ ，并且 $\mathbf{g}(\mathbf{x} | \theta)$ 是根据参数 θ 的集合定义的从 d 维到 k 维的映射函数。前面我们讨论的经典的 MDS 对应线性变换

$$\mathbf{z} = \mathbf{g}(\mathbf{x} | \mathbf{W}) = \mathbf{W}^T \mathbf{x} \quad (6.28)$$

但是在一般情况下,也可以使用非线性的映射;这称作 Sammon 映射(Sammon mapping)。在映射中的标准化误差称作 Sammon 应力(Sammon stress),定义为

$$\begin{aligned} E(\theta | \mathcal{X}) &= \sum_{r,s} \frac{(\|\mathbf{z}^r - \mathbf{z}^s\| - \|\mathbf{x}^r - \mathbf{x}^s\|)^2}{\|\mathbf{x}^r - \mathbf{x}^s\|^2} \\ &= \sum_{r,s} \frac{(\|\mathbf{g}(\mathbf{x}^r | \theta) - \mathbf{g}(\mathbf{x}^s | \theta)\| - \|\mathbf{x}^r - \mathbf{x}^s\|)^2}{\|\mathbf{x}^r - \mathbf{x}^s\|^2} \end{aligned} \quad (6.29)$$

可以对 $\mathbf{g}(\cdot | \theta)$ 使用任何回归方法,估计 θ 来最小化训练数据 \mathcal{X} 上的应力。如果 $\mathbf{g}(\cdot)$ 在 \mathbf{x} 上是非线性的,这将对一个非线性的维度归约。

在分类的情况下,我们可以在距离中包含类信息(见 Webb, 1999),如

$$d'_{rs} = (1 - \alpha)d_{rs} + \alpha c_{rs}$$

其中 c_{rs} 是 \mathbf{x}^r 和 \mathbf{x}^s 所属类之间的“距离”。这个类间距离应该被主观地提供,而 α 用交叉确认优化。

6.6 线性判别分析

线性判别分析(linear discriminant analysis, LDA)是一种用于分类问题的维度归约的监督方法。我们由两个类来开始这个问题的讨论,然后推广到 $K > 2$ 个类。

给定来自两个类 C_1 和 C_2 的样本,我们希望找到由向量 \mathbf{w} 定义的方向,使得当数据投影到 \mathbf{w} 上时,来自两个类的样本尽可能地分开。正如我们在前面看到的

$$\mathbf{z} = \mathbf{w}^T \mathbf{x} \quad (6.30)$$

是 \mathbf{x} 到 \mathbf{w} 上的投影,因而也是一个从 d 维到 1 维的维度归约。

128

\mathbf{m}_1 和 m_1 分别是 C_1 类样本在投影前和投影后的均值。注意 $\mathbf{m}_1 \in \mathbb{R}^d$ 而 $m_1 \in \mathbb{R}$ 。我们有样本 $\mathcal{X} = \{\mathbf{x}^i, r^i\}$, 使得如果 $\mathbf{x}^i \in C_1$ 则 $r^i = 1$, 而如果 $\mathbf{x}^i \in C_2$ 则 $r^i = 0$ 。

$$\begin{aligned} m_1 &= \frac{\sum_i \mathbf{w}^T \mathbf{x}^i r^i}{\sum_i r^i} = \mathbf{w}^T \mathbf{m}_1 \\ m_2 &= \frac{\sum_i \mathbf{w}^T \mathbf{x}^i (1 - r^i)}{\sum_i (1 - r^i)} = \mathbf{w}^T \mathbf{m}_2 \end{aligned} \quad (6.31)$$

来自 C_1 和 C_2 的样本投影后的散布(scatter)是

$$\begin{aligned} s_1^2 &= \sum_i (\mathbf{w}^T \mathbf{x}^i - m_1)^2 r^i \\ s_2^2 &= \sum_i (\mathbf{w}^T \mathbf{x}^i - m_2)^2 (1 - r^i) \end{aligned} \quad (6.32)$$

投影后,为了使两个类被很好地分开,我们希望均值尽可能远离,并且类实例散布在尽可能小的区域中。因此,我们希望 $|m_1 - m_2|$ 大,而 $s_1^2 + s_2^2$ 小(见图 6-7)。费希尔线性判别式(Fisher's linear discriminant)是最小化式(6.33)的 \mathbf{w} 。

$$J(\mathbf{w}) = \frac{(m_1 - m_2)^2}{s_1^2 + s_2^2} \quad (6.33)$$

129

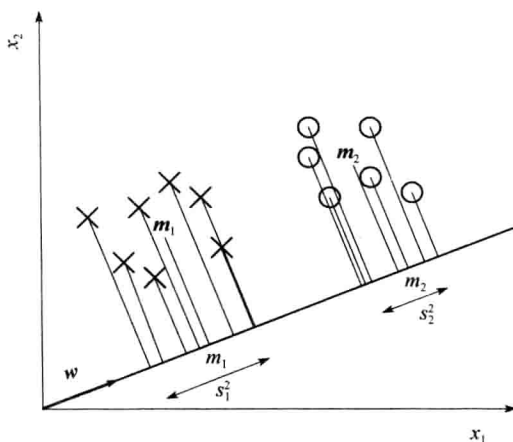


图 6-7 两维、两类的数据在 \$\mathbf{w}\$ 上的投影

重写分子，我们得到

$$\begin{aligned}
 (\mathbf{m}_1 - \mathbf{m}_2)^2 &= (\mathbf{w}^T \mathbf{m}_1 - \mathbf{w}^T \mathbf{m}_2)^2 \\
 &= \mathbf{w}^T (\mathbf{m}_1 - \mathbf{m}_2) (\mathbf{m}_1 - \mathbf{m}_2)^T \mathbf{w} \\
 &= \mathbf{w}^T \mathbf{S}_B \mathbf{w}
 \end{aligned} \tag{6.34}$$

其中 \$\mathbf{S}_B = (\mathbf{m}_1 - \mathbf{m}_2) (\mathbf{m}_1 - \mathbf{m}_2)^T\$ 是类间散布矩阵 (between-class scatter matrix)。这个分子是投影后类实例在其均值周围散布的和，并且可以改写为

$$\begin{aligned}
 s_1^2 &= \sum_i (\mathbf{w}^T \mathbf{x}^i - m_1)^2 r^i \\
 &= \sum_i \mathbf{w}^T (\mathbf{x}^i - \mathbf{m}_1) (\mathbf{x}^i - \mathbf{m}_1)^T \mathbf{w} r^i \\
 &= \mathbf{w}^T \mathbf{S}_1 \mathbf{w}
 \end{aligned} \tag{6.35}$$

其中

$$\mathbf{S}_1 = \sum_i r^i (\mathbf{x}^i - \mathbf{m}_1) (\mathbf{x}^i - \mathbf{m}_1)^T \tag{6.36}$$

是 \$C_1\$ 的类内散布矩阵 (within-class scatter matrix)。\$\mathbf{S}_1 / \sum_i r^i\$ 是 \$\mathbf{\Sigma}_1\$ 的估计。类似地，\$s_2^2 = \mathbf{w}^T \mathbf{S}_2 \mathbf{w}\$，其中 \$\mathbf{S}_2 = \sum_i (1 - r^i) (\mathbf{x}^i - \mathbf{m}_2) (\mathbf{x}^i - \mathbf{m}_2)^T\$，并且我们得到

$$s_1^2 + s_2^2 = \mathbf{w}^T \mathbf{S}_w \mathbf{w}$$

其中，\$\mathbf{S}_w = \mathbf{S}_1 + \mathbf{S}_2\$ 是类内散布的总和。注意，\$s_1^2 + s_2^2\$ 除以样本总数是汇聚数据的方差。式 (6.33) 可以改写为

$$J(\mathbf{w}) = \frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}{\mathbf{w}^T \mathbf{S}_w \mathbf{w}} = \frac{|\mathbf{w}^T (\mathbf{m}_1 - \mathbf{m}_2)|^2}{\mathbf{w}^T \mathbf{S}_w \mathbf{w}} \tag{6.37}$$

关于 \$\mathbf{w}\$ 取 \$J\$ 的导数并令其等于 0，我们得到

$$\frac{\mathbf{w}^T (\mathbf{m}_1 - \mathbf{m}_2)}{\mathbf{w}^T \mathbf{S}_w \mathbf{w}} \left(2(\mathbf{m}_1 - \mathbf{m}_2) - \frac{\mathbf{w}^T (\mathbf{m}_1 - \mathbf{m}_2)}{\mathbf{w}^T \mathbf{S}_w \mathbf{w}} \mathbf{S}_w \mathbf{w} \right) = 0$$

给定 \$\mathbf{w}^T (\mathbf{m}_1 - \mathbf{m}_2) / \mathbf{w}^T \mathbf{S}_w \mathbf{w}\$ 为常数，我们有

$$\mathbf{w} = c\mathbf{S}_w^{-1}(\mathbf{m}_1 - \mathbf{m}_2) \quad (6.38)$$

其中 c 是某个常数。因为对我们来说重要的是方向，而不是大小，所以我们可以取 $c = 1$ 并找出 \mathbf{w} 。

记住当 $p(\mathbf{x} | C_i) \sim \mathcal{N}(\boldsymbol{\mu}_i, \boldsymbol{\Sigma})$ 时，我们有线性判别式，其中 $\mathbf{w} = \boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)$ ，并且我们看到如果类是正态分布的，则费希尔线性判别式是最优的。在同样的假设下，我们还可以计算阈值 w_0 ，来分开两个类。但是，费希尔线性判别式甚至在类不是正态分布时也能使用。我们已经把样本从 d 维投影到 1 维，之后可以使用任何分类方法。

在 $K > 2$ 个类的情况下，我们希望找到这样的矩阵 \mathbf{W} ，使得

$$\mathbf{z} = \mathbf{W}^T \mathbf{x} \quad (6.39)$$

其中 \mathbf{z} 是 k 维的， \mathbf{W} 是 $d \times k$ 矩阵。 C_i 的类内散布矩阵是

$$\mathbf{S}_i = \sum_t r'_i(\mathbf{x}^t - \mathbf{m}_i)(\mathbf{x}^t - \mathbf{m}_i)^T \quad (6.40)$$

其中，如果 $\mathbf{x}^t \in C_i$ 则 $r'_i = 1$ ，否则为 0。总类内散布是

$$\mathbf{S}_w = \sum_{i=1}^K \mathbf{S}_i \quad (6.41)$$

当存在 $K > 2$ 个类时，均值的散布根据它们在总均值周围的散布情况计算

$$\mathbf{m} = \frac{1}{K} \sum_{i=1}^K \mathbf{m}_i \quad (6.42)$$

而类间散布矩阵是

$$\mathbf{S}_B = \sum_{i=1}^K N_i(\mathbf{m}_i - \mathbf{m})(\mathbf{m}_i - \mathbf{m})^T \quad (6.43)$$

其中 $N_i = \sum_t r'_i$ 。投影后的类间散布矩阵是 $\mathbf{W}^T \mathbf{S}_B \mathbf{W}$ ，而投影后的类内散布矩阵是 $\mathbf{W}^T \mathbf{S}_w \mathbf{W}$ 。它们都是 $k \times k$ 矩阵。我们希望第一个散布大；也就是说，在投影之后，在新的 k 维空间，我们希望类均值互相之间尽可能远离。我们希望第二个散布小；也就是说，在投影之后，我们希望来自同一个类的样本尽可能接近它们的均值。对于一个散布(或协方差)矩阵，散布的一个度量是行列式。记住该行列式是特征值的乘积，而特征值给出沿着它的特征向量(成分)的方差。因此，我们对最小化式(6.44)的矩阵 \mathbf{W} 感兴趣

$$J(\mathbf{W}) = \frac{|\mathbf{W}^T \mathbf{S}_B \mathbf{W}|}{|\mathbf{W}^T \mathbf{S}_w \mathbf{W}|} \quad (6.44)$$

$\mathbf{S}_w^{-1} \mathbf{S}_B$ 的最大的特征向量是解。 \mathbf{S}_B 是 K 个秩为 1 的矩阵 $(\mathbf{m}_i - \mathbf{m})(\mathbf{m}_i - \mathbf{m})^T$ 的和，并且它们之中只有 $K - 1$ 个是独立的。因此， \mathbf{S}_B 具有最大秩 $K - 1$ ，并且我们取 $k = K - 1$ 。这样，我们定义一个新的、较低的 $K - 1$ 维空间，然后在那里构造判别式(见图 6-8)。虽然 LDA 使用类分离性作为它的好坏标准，但是在这个新空间里可以使用任意的分类方法来估计判别式。

我们看到，为了使用 LDA， \mathbf{S}_w 应该是可逆的。如果不可逆，则我们可以先用 PCA 消除奇异性，然后把 LDA 应用于其结果。然而，我们应该确保 PCA 不会把维度降低太多，从而使得 LDA 没有多少事可做。

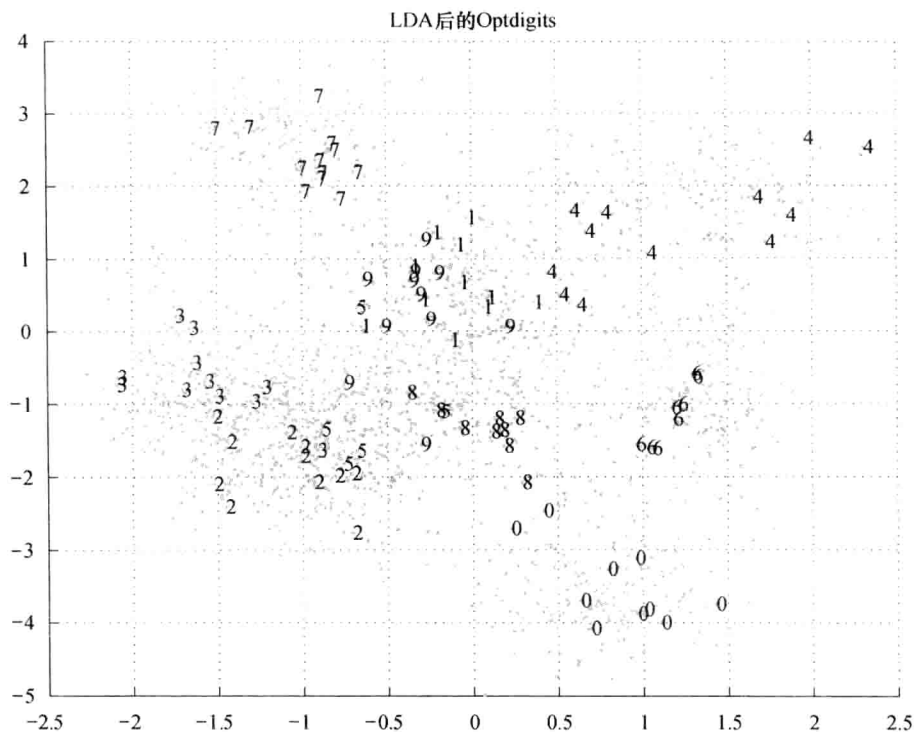


图 6-8 绘制在 LDA 找到的前两个维空间上的 Optdigits。与图 6-3 比较，正如期望的那样，我们看到 LDA 比 PCA 导致更好的类分离。即便在这个二维空间(有 9 个)，我们也能看到不同类的分开的云团

6.7 等距特征映射

当数据落在一个线性子空间中时，6.3 节讨论的主成分分析 (PCA) 效果很好。然而，在许多应用中，这一前提并不成立。以人脸识别为例，人脸用如 100×100 的二维图像表示，在这种情况下，每张人脸是 10 000 维空间中的一个点。现在，假设随着一个人由右向左慢慢转动头部，我们取一系列照片。我们得到的这个面部图像序列沿着 10 000 维空间中的一条轨迹，而这条曲线不是线性的。现在，考虑许多人的脸，随着他们转动头部，所有人的脸轨迹定义了 10 000 维空间中的一个流形，并且这就是我们想要建模的。两张人脸的相似性不能简单地表示成像素差的和，因而欧氏距离不是一个好的度量。很可能出现这种情况：两个具有相同姿势的不同人的图像的欧氏距离比同一个人的两种不同姿势图像的欧氏距离更小，这不是我们想要的。我们应该计算的是沿流形的距离，这称作测地距离 (geodesic distance)。等距特征映射 (Isometric feature mapping, Isomap) (Tenenbaum、de Silva 和 Langford 2000) 估计这种距离，并使用多维定标 (multidimensional scaling, MDS) (见 6.5 节)，用它进行维度归约。

Isomap 使用所有数据点对之间的测地距离。对于输入空间中靠近的邻近点，可以使用欧氏距离；对于姿势中的小改变，流形是局部线性的。对于远离的点，测地距离用沿流形的

点之间的距离和来近似。可以这样做：定义一张图，其结点对应于 N 个数据点，其边连接邻近的点（距离小于某个 ε 的点，或 n 个最近邻之一），边的权重对应于欧氏距离。任意两点之间的测地距离用对应的两点之间最短路径长度计算。对于两个不邻近的点，我们需要沿通路跳过许多中间点，因而该距离是沿流形的距离，用局部欧氏距离的和来近似（参见图 6-9）。

如果 $\|\mathbf{x}^r - \mathbf{x}^s\| < \varepsilon$ （同时要确保图是连通的）或者 \mathbf{x}^s 是 \mathbf{x}^r 的 n 个最近邻之一（同时要确保距离矩阵是对称的），并且设置其边长为 $\|\mathbf{x}^r - \mathbf{x}^s\|$ ，则两个结点 r 和 s 是连接的。对于两个结点 r 和 s ， d_{rs} 是它们之间最短路径的长度。然后，在 d_{rs} 上应用 MDS，通过观察它们贡献的方差比例，把维度降低到 k 。这具有如下效果：把测地空间中相隔较远的 r 和 s 也放在新的 k 维空间相隔较远的位置上，即使在原 d 维空间中它们在欧氏距离意义下靠近也如此。

很显然，随着点数增加，图形距离提供了更好的近似，尽管以更多的运行时间为代价。如果时间是至关重要的，则可以二次抽样并使用“地标点”子集使得算法更快。参数 ε 需要细心调整，如果它太小，则可能存在多个连通分支；而如果它太大，则可能添加破坏低维嵌入的“捷径”边（Balasubramanian 等 2002）。

正如使用 MDS 一样，使用 Isomap 的一个问题是：它把 N 个点放到一个低维空间，但是它并不学习一个可以映射新检验点的一般映射函数；新点应该添加到数据集中，并且需要使用 $N+1$ 个实例再次运行整个算法。

6.8 局部线性嵌入

局部线性嵌入（locally linear embedding, LLE）由局部线性拟合发现全局非线性结构（Roweis 和 Saul 2000）。其基本思想是，流形的每个局部小段都可以线性地近似，并且给定足够多的数据，每个点都可以表示成其近邻（或者用给定的近邻数 n ，或者用距离阈值 ε 定义）的线性加权和。给定原空间中的 \mathbf{x}^r 和它的近邻 $\mathbf{x}_{(r)}^s$ ，我们可以使用最小二乘找到重构权重 \mathbf{W}_n 以最小化误差函数

$$\mathcal{E}^w(\mathbf{W}|\mathcal{X}) = \sum_r \left\| \mathbf{x}^r - \sum_s \mathbf{W}_{rs} \mathbf{x}_{(r)}^s \right\|^2 \quad (6.45)$$

约束条件是：对于任意 r ， $\sum_s \mathbf{W}_{rs} = 1$ ，并且 $\sum_s \mathbf{W}_{rs} = 0$ 。

LLE 的基本思想是，重构权重 \mathbf{W}_{rs} 以反映数据的固有几何性质，我们期望这种性质对于流形（即实例映射到的新空间）的局部小段也有效（见图 6.10）。因此，LLE 的第二步是保持权重 \mathbf{W}_n 固定，并令新坐标 \mathbf{z}^r 取满足由权重给定的内插点约束所需的值：

$$\mathcal{E}^z(\mathbf{Z}|\mathbf{W}) = \sum_r \left\| \mathbf{z}^r - \sum_s \mathbf{W}_{rs} \mathbf{z}^s \right\|^2 \quad (6.46)$$

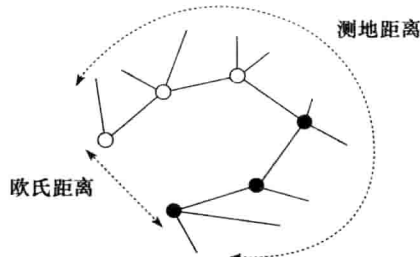


图 6-9 与欧氏距离不同，测地距离沿流形计算，而欧氏距离不使用这种信息。多维定标之后，来自两个类的两个实例将被映射到新空间中相隔较远的位置，尽管它们在原空间很靠近

在原 d 维空间位于附近的点在新的 k 维空间应该保持邻近, 并且类似地新的 k 维空间中协同定位。式(6.46)可以改写为

$$\mathcal{E}^z(\mathbf{Z} | \mathbf{W}) = \sum_{r,s} \mathbf{M}_{rs} (\mathbf{z}^r)^T \mathbf{z}^s \quad (6.47)$$

其中

$$\mathbf{M}_{rs} = \delta_{rs} - \mathbf{W}_{rs} - \mathbf{W}_{sr} + \sum_i \mathbf{W}_{ir} \mathbf{W}_{is} \quad (6.48)$$

\mathbf{M} 是稀疏的(一个数据点的近邻只占很小的比例: $n \ll N$)、对称的和半正定的。与其他维归约方法一样, 我们要求数据在中心点, 即 $E[\mathbf{z}] = 0$, 并且新坐标是不相关的且具有单位长度: $\text{Cov}(\mathbf{z}) = \mathbf{I}$ 。在这两个约束下, 式(6.47)的解由具有最小特征值的 $k+1$ 个特征向量给出。我们忽略特征值中最小的一个, 其余 k 个特征向量即为新坐标。

135

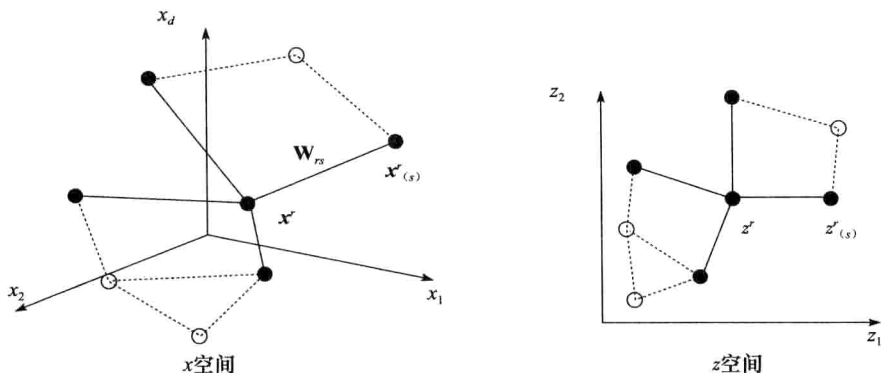


图 6-10 局部线性嵌入首先在原空间学习约束, 再关于根据约束把点放置在新空间。约束使用直接近邻(用实线显示)学习, 但是也传播到二级近邻(用虚线显示)

因为 n 个近邻生成一个 $n-1$ 维空间(在二维空间中, 需要到三个点的距离来唯一确定你的位置), 因此 LLE 可以把维度归约到 $k \leq n-1$ 。据观察(Saul 和 Roweis 2003), k 与 n 之间的某个范围足以得到一个好的嵌入。注意, 如果 n (或 ε) 很小, 则通过把每个实例与它的近邻连接构造的图可能不再连通, 并且可能需要在各连通分支上分别运行 LLE, 以找出输入空间不同部分的流形。另一方面, 如果 n (或 ε) 取值太大, 那么某些近邻可能太远, 以使局部线性假设不成立, 可能损害嵌入。有可能基于某种先验知识, 在输入空间的不同部分使用不同的 n (或 ε), 但是如何做仍然是一个尚需研究的问题(Saul 和 Roweis 2003)。

与 Isomap 一样, LLE 解是 N 个点的新坐标系, 但我们不学习映射, 因而不能针对新的 \mathbf{x}' 找到 \mathbf{z}' 。对此有两种解决方案:

1) 使用相同的思想, 可以在原 d 维空间找出 \mathbf{x}' 的 n 个近邻, 并且首先学习最小化

$$\mathcal{E}^w(\mathbf{w} | \mathcal{X}) = \left\| \mathbf{x}' - \sum_s \mathbf{w}_s \mathbf{x}^s \right\|^2 \quad (6.49)$$

的重构权重 \mathbf{w}_j , 然后使用它们在新的 k 维空间中重构 \mathbf{z}' :

$$\mathbf{z}' = \sum_s \mathbf{w}_s \mathbf{z}^s \quad (6.50)$$

注意：这种方法也可以用于由 Isomap(或 MDS)解插值。然而，其缺点是需要存放整个数据集 $\{\mathbf{x}', \mathbf{z}'\}_{i=1}^N$ 。

2) 使用 $\mathcal{X} = \{\mathbf{x}', \mathbf{z}'\}_{i=1}^N$ 作为训练集，可以训练任意回归器 $\mathbf{g}(\mathbf{x}' | \theta)$ (例如，一个多层感知器(第 11 章))，作为由 \mathbf{x}' 近似 \mathbf{z}' 的泛化器，其参数 θ 通过学习得到，以最小化回归误差：

$$\mathcal{E}(\theta | \mathcal{X}) = \sum_i \|\mathbf{z}' - \mathbf{g}(\mathbf{x}' | \theta)\|^2 \quad (6.51)$$

一旦训练完成，我们就可以计算 $\mathbf{z}' = \mathbf{g}(\mathbf{x}' | \theta)$ 。模型 $\mathbf{g}(\cdot)$ 应该小心选择，使之能够学习映射。可能不再是唯一最优的，因而存在通常与最小化有关的所有问题，即初始化、局部最优、收敛等。

在 Isomap 和 LLE 中，存在局部信息，该信息在近邻上传播，以得到全局解。在 Isomap 中，测地距离是局部距离的和；在 LLE 中，最终放置 \mathbf{z}' 的最优解考虑了所有局部 \mathbf{W}_n 值。假设 a 和 b 是近邻， b 和 c 是近邻。尽管 a 和 c 可能不是近邻，但是 a 和 c 之间的依赖性或者通过图， $d_{ac} = d_{ab} + d_{bc}$ ，或者通过权重 \mathbf{W}_{ab} 和 \mathbf{W}_{bc} 而存在。在这两个算法中，全局非线性组织通过整合部分重叠的局部线性约束而发现。

6.9 注释

特征选择算法的一个综述在 Devijer 和 Kittler 1982 中给出。特征子集选择算法又称包装(wrapper)方法，其中特征选择被看作“包裹”在学习方法上作为子例程使用(Kohavi 和 John 1997)。Miller 1990 讨论了回归中的子集选择。我们讨论的向前和向后搜索过程是局部搜索过程。Fukunaga 和 Narendra(1977)提出了一种分支和限制(branch 和 bound)过程。以更大的开销为代价，我们可以使用诸如模拟退火或者遗传算法这样的随机过程，在搜索空间中进行更广泛的搜索。

还有一些用于特征选择的过滤(filtering)算法，其中启发式度量在预处理阶段用来计算特征的“相关性”，而不实际使用学习方法。例如，对于分类，取代每步训练和检验一个分类器，我们可以使用像在线性判别式分析所用的可分性度量来度量在新空间中把类彼此分开的质量(McLachlan 1992)。随着计算费用的降低，最好在循环中包含学习方法，因为不能保证过滤方法使用的启发式度量与使用特征的学习方法的偏倚匹配。启发式度量都不能取代实际的验证准确率。Guyon 和 Elisseeff(2003)给出了特征选择方法的综述。

投影方法需要数值输入，并且离散变量应该用 0/1 哑变量表示，而子集选择可以直接使用离散输入。寻找特征向量和特征值是非常直接的，一个代码的例子在 Press 等 1992 中给出。因子分析是由英国的心理学家 Charles Spearman 引入的，用于发现智力的单个因素，来解释不同智力测试得分之间的联系。这种单个因子(称为 g)的存在性是讨论的课题。更多关于多维定标的信息可以在 Cox 和 Cox 1994 中找到。

我们讨论的投影方法是批过程，因为它们要求在发现投影方向之前给定整个的样本。Mao 和 Jain(1995)讨论做 PCA 和 LDA 的在线过程，其中样例被逐个给出，并且更新随新实例的到达而进行。做非线性投影的另外一种可能的方法是 Sammon 映射的估计被取作非线性

136
137

函数,例如,多层感知器(见 11.11 节)(Mao 和 Jain 1995)。进行非线性的因子分析也是可能的,但是相当困难。当模型是非线性的时,构建一个正确的非线性模型是很困难的。我们还需要用复杂的优化和逼近方法来求解模型参数。

更多的信息可以参阅 Isomap 的主页 <http://web.mit.edu/cocosci/isomap.html> 和 LLE 的主页 <http://www.cs.toronto.edu/~roweis/lle/>, 它们都包含相关发表物和实例代码的链接。

正如我们通过把高阶项看作附加的输入,使用线性回归来实现多项式回归(见 5.8 节)一样,另一种实现非线性维归约的方法是首先使用非线性基函数映射到新空间,然后在新空间使用线性模型。在第 13 章我们讨论核方法将会看到如何有效地实现这一做法。

特征提取和决策制定之间有一个权衡。如果特征提取很好,分类(或回归)的任务变得微不足道。例如,当类代码被作为新的特征从现有特征中提取出来时就是如此。另一方面,如果分类方法足够好,则没有必要进行特征提取;它自己做它的自动特征选择或者内部组合。我们处于这两个理想世界之间。

存在一些算法,它们内部做某些特征选择,尽管是以有限的方式。决策树(第 9 章)在产生决策树时进行特征选择,而多层感知器(第 11 章)在隐藏节点做非线性的特征提取。我们期望沿着这个方向,在耦合特征提取和其后的分类或回归方面取得更多进展。

6.10 习题

1. 假定类是正态分布的,在子集选择中,当一个变量被添加或去掉时,如何快速计算新的判别式?例如,如何从 S_{old}^{-1} 计算 S_{new}^{-1} ?
2. 使用取自 UCI 知识库中的 Optdigits 实现 PCA。对于不同数量的特征向量,重构数字图像并计算重构误差式(6.12)。
3. 给定道路旅行距离作为输入,使用 MDS 绘制你所在州或国家的地图。
4. 在 Sammon 映射中,如果映射是线性的,即 $g(\mathbf{x} | \mathbf{W}) = \mathbf{W}^T \mathbf{x}$,如何计算最小化 Sammon 应力的 \mathbf{W} ?
5. 重做习题 3。这次使用 Isomap,其中仅当两个城市之间有不经过其他城市的直通道路时,才连接这两个城市。
6. 在 Isomap 中,我们也可以使用邻近点之间的马氏距离,而不使用欧氏距离。如果有的话,这种方法有何优缺点?
7. 绘制两类、二维数据,使得(a)PCA 和 LDA 找到相同的方向;(b)PCA 和 LDA 找到完全不同的方向。
8. 只要我们有对象两两之间的距离,就可以使用多维定标。只要我们有某种相似性度量,就完全不必把对象用向量表示。你能给出一个例子吗?
9. 如何把类信息结合到 Isomap 和 LLE 中,使得相同类的距离被映射到新空间的邻近位置?
10. 在因子分析中,如果我们已经知道某些因子,如何找到其余的因子?
11. 讨论一个应用,其中有隐藏因子(不必是线性的),并且因子分析有望效果很好。

138
139

6.11 参考文献

- Balasubramanian, M., E. L. Schwartz, J. B. Tenenbaum, V. de Silva, and J. C. Langford. 2002. "The Isomap Algorithm and Topological Stability." *Science* 295: 7.
- Chatfield, C., and A. J. Collins. 1980. *Introduction to Multivariate Analysis*. London: Chapman and Hall.
- Cox, T. F., and M. A. A. Cox. 1994. *Multidimensional Scaling*. London: Chapman and Hall.
- Devijver, P. A., and J. Kittler. 1982. *Pattern Recognition: A Statistical Approach*. New York: Prentice-Hall.
- Flury, B. 1988. *Common Principal Components and Related Multivariate Models*. New York: Wiley.
- Fukunaga, K., and P. M. Narendra. 1977. "A Branch and Bound Algorithm for Feature Subset Selection." *IEEE Transactions on Computers* C-26: 917-922.
- Guyon, I., and A. Elisseeff. 2003. "An Introduction to Variable and Feature Selection." *Journal of Machine Learning Research* 3: 1157-1182.
- Hastie, T. J., R. J. Tibshirani, and A. Buja. 1994. "Flexible Discriminant Analysis by Optimal Scoring." *Journal of the American Statistical Association* 89: 1255-1270.
- Kohavi, R., and G. John. 1997. "Wrappers for Feature Subset Selection." *Artificial Intelligence* 97: 273-324.
- Mao, J., and A. K. Jain. 1995. "Artificial Neural Networks for Feature Extraction and Multivariate Data Projection." *IEEE Transactions on Neural Networks* 6: 296-317.
- McLachlan, G. J. 1992. *Discriminant Analysis and Statistical Pattern Recognition*. New York: Wiley.
- Miller, A. J. 1990. *Subset Selection in Regression*. London: Chapman and Hall.
- Press, W. H., B. P. Flannery, S. A. Teukolsky, and W. T. Vetterling. 1992. *Numerical Recipes in C*. Cambridge, UK: Cambridge University Press.
- Pudil, P., J. Novovičová, and J. Kittler. 1994. "Floating Search Methods in Feature Selection." *Pattern Recognition Letters* 15: 1119-1125.
- Rencher, A. C. 1995. *Methods of Multivariate Analysis*. New York: Wiley.
- Roweis, S. T., and L. K. Saul. 2000. "Nonlinear Dimensionality Reduction by Locally Linear Embedding." *Science* 290: 2323-2326.
- Saul, K. K., and S. T. Roweis. 2003. "Think Globally, Fit Locally: Unsupervised Learning of Low Dimensional Manifolds." *Journal of Machine Learning Research* 4: 119-155.
- Tenenbaum, J. B., V. de Silva, and J. C. Langford. 2000. "A Global Geometric Framework for Nonlinear Dimensionality Reduction." *Science* 290: 2319-2323.
- Tipping, M. E., and C. M. Bishop. 1999. "Probabilistic Principal Components Analysis." *Journal of the Royal Statistical Society Series B* 61: 611-622.
- Turk, M., and A. Pentland. 1991. "Eigenfaces for Recognition." *Journal of Cognitive Neuroscience* 3: 71-86.
- Webb, A. 1999. *Statistical Pattern Recognition*. London: Arnold.

140

141

聚 类

在参数方法中，我们假设样本来自一个已知的分布。当这种假设站不住脚时，我们放宽该假设，并使用半参数方法，允许用混合分布估计输入样本。聚类方法允许从数据中学习混合分布。除了概率建模之外，我们还讨论向量量化和层次聚类。

7.1 引言

在第4章和第5章中，我们讨论了密度估计的参数方法，在那里我们假设样本 \mathbf{x} 取自某个参数族，例如高斯族。在参数分类中，这对应为类密度 $p(\mathbf{x} | C_i)$ 假定某种密度。参数方法的优点是，给定一个模型，问题归结为少量参数的估计。对于密度估计，参数是密度的充分统计量。例如，对于高斯密度，参数为均值和协方差。

尽管参数方法使用频繁，但是对于假定并不成立的许多应用来说，假定一个严格的参数模型可能是偏倚根源。因此我们需要更灵活的模型。尤其是，假定高斯密度对应假定样本（例如一个类的实例）形成 d -维空间中的单个分组，并且正如我们在第5章所看到的，该分组的中心和形状分别由均值协方差给定。

然而，在许多应用中，样本不是一个分组，而可能有多个分组。以手写字符识别为例：有两种风格书写数字7；美洲人的写法是‘7’，而欧洲人的写法是中间有一个水平杠（与欧洲人手写的、上面有一小划的‘1’以示区别）。在这种情况下，当样本包含来自两个洲的实例时，数字7应当表示成两个不相交的分组。如果每个分组用一个高斯分布表示，则该类可以用两个高斯分布的混合分布表示，每个高斯分布用于一种书写风格。

类似的例子是语音识别，其中由于不同的发音、口音、性别、年龄等，相同的词可能以不同的方法读出。这样，当没有单个、通用的原型时，为了统计上的正确性，应当在密度中表示所有这些不同的方法。

我们称这种方法为半参数密度估计(semiparametric density estimation)，因为我们仍然对样本中的每个分组假定一个参数模型。在第8章，我们将讨论非参数方法。当数据没有结构，甚至连混合模型都不能使用时，可以使用非参数方法。本章我们致力于密度估计，而将监督学习放到第12章。

7.2 混合密度

混合密度(mixture density)记作

$$p(\mathbf{x}) = \sum_{i=1}^k p(\mathbf{x} | G_i) P(G_i) \quad (7.1)$$

其中 G_i 是混合分支(mixture component)，也称分组(group)或簇(cluster)。 $p(\mathbf{x} | G_i)$ 是支密度

142
143

(component density), 而 $P(G_i)$ 是混合比例 (mixture proportion)。分支数 k 是超级参数, 应当预先指定。给定样本和 k , 学习对应估计支密度和比例。当我们假定支密度遵守参数模型时, 我们只需要估计它们的参数。如果支密度是多元高斯的, 则我们有 $p(\mathbf{x} | G_i) \sim \mathcal{N}(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$, 而 $\Phi = \{P(G_i), \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i, \}_{i=1}^k$ 是应当从独立同分布的样本 $\mathcal{X} = \{\mathbf{x}^t\}_t$ 中学习的参数。

[144]

参数分类是名副其实的混合模型, 其中分组 G_i 对应类 C_i , 支密度 $p(\mathbf{x} | G_i)$ 对应类密度 $p(\mathbf{x} | C_i)$, 而 $P(G_i)$ 对应类先验 $P(C_i)$:

$$p(\mathbf{x}) = \sum_{i=1}^K p(\mathbf{x} | C_i) P(C_i)$$

在这种监督情况下, 我们知道有多少个分组, 而学习参数却是平凡的, 因为我们有类标号, 即知道哪个实例属于哪个类 (分支)。从第 5 章我们知道, 给定样本 $\mathcal{X} = \{\mathbf{x}^t, \mathbf{r}^t\}_{t=1}^N$, 其中如果 $\mathbf{x}^t \in C_i$ 则 $r_i^t = 1$, 否则 r_i^t 为 0, 可以使用最大似然计算这些参数。当每个类都是高斯分布时, 我们有混合高斯分布, 并且参数估计为

$$\begin{aligned}\hat{P}(C_i) &= \frac{\sum_t r_i^t}{N} \\ \mathbf{m}_i &= \frac{\sum_t r_i^t \mathbf{x}^t}{\sum_t r_i^t} \\ S_i &= \frac{\sum_t r_i^t (\mathbf{x}^t - \mathbf{m}_i)(\mathbf{x}^t - \mathbf{m}_i)^T}{\sum_t r_i^t}\end{aligned}\quad (7.2)$$

本章的不同是样本为 $\mathcal{X} = \{\mathbf{x}^t\}_t$, 我们有非监督学习 (unsupervised learning) 问题。我们只有 \mathbf{x}^t 而没有标号 \mathbf{r}^t , 也就是说我们不知道 \mathbf{x}^t 来自哪个分支。这样, 我们应当估计二者: 第一, 我们应当估计给定实例所属的分支标号 r_i^t ; 第二, 一旦我们估计了标号, 我们就要估计给定实例集所属分支的参数。为此, 我们首先讨论一种简单的聚类算法 k -均值, 并在后面证明它是期望最大化 (Expectation-Maximization) 算法的一个特例。

7.3 k -均值聚类

假设我们有一幅图像, 按 24 位/像素存放, 而我们可能有多达 1 600 万种颜色。假定我们有 8 位/像素的彩色屏幕, 只能显示 256 种颜色。我们想在 1 600 万种颜色中找出最佳的 256 种颜色, 使得仅用了调色板中 256 种颜色的图像看上去尽可能接近原来的图像。这是颜色量化 (color quantization) 问题, 其中我们从高分辨率映射到低分辨率。在一般情况下, 目标是从连续空间映射到离散空间; 这一过程称作向量量化 (vector quantization)。

[145]

当然, 我们总能均匀地进行量化, 但是把映射表指派到图像中不存在的颜色, 或不给图像中频繁使用的颜色分配附加的表目会错失颜色映射。例如, 如果图像是海景, 则我们可望看到许多深浅不一的蓝色而可能不是红色。因此, 颜色映射表目的分布应当尽可能接近地反映原来的密度, 将更多的表目放在高密度区域, 而丢弃没有数据的区域。

假定我们有样本 $\mathcal{X} = \{\mathbf{x}^t\}_{t=1}^N$ 。我们有 k 个参考向量 (reference vector) $\mathbf{m}_j, j = 1, \dots, k$ 。在我们的颜色量化的例子中, \mathbf{x}^t 是 24 位的图像像素值, \mathbf{m}_j 是颜色映射表目, 也是 24 位,

$k = 256$ 。

暂时假定我们以某种方法得到了 m_j 的值；稍后我们将讨论如何学习它们。为了显示图像，给定像素 x' ，我们用颜色映射中最相似的、满足下式的表目 m_i 表示它

$$\|x' - m_i\| = \min_j \|x' - m_j\|$$

也就是说，我们使用参考向量符号系统中最接近的值，而不是使用原始数据。 m_i 又称码本向量 (codebook vector) 或码字 (code word)，因为这是一个编码/解码过程 (参见图7-1)：从 x' 到 i 是使用编码本 $m_i (i=1, \dots, k)$ 对数据编码的过程，而在接收端，从 i 产生 m_i 是解码。量化也能压缩：例如，替代使用 24 位存储 (或在通讯线上传输) 每个 x' ，我们可以只存储/传输它在颜色映射中的下标 i ，使用 8 位索引 1 到 256 中的值，我们得到几乎为 3 的压缩率；存储/传输的也是颜色映射。

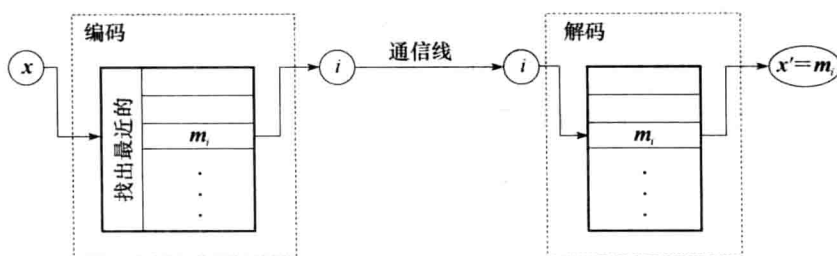


图 7-1 给定 x ，编码器发送最近的码字的下标，而解码器使用接收到的下标产生码字 x' 。误差为 $\|x' - x\|^2$

让我们看看如何计算 m_i ：当 x' 用 m_i 表示时，存在一个正比于距离 $\|x' - m_i\|$ 的误差。为了使新图像看上去像原来的图像，我们应当针对所有的像素，使该距离尽可能小。总重构误差 (reconstruction error) 定义为

$$E(\{m_i\}_{i=1}^k | X) = \sum_i \sum_t b_i^t \|x^t - m_i\|^2 \quad (7.3)$$

其中

$$b_i^t = \begin{cases} 1 & \text{如果 } \|x^t - m_i\| = \min_j \|x^t - m_j\| \\ 0 & \text{否则} \end{cases} \quad (7.4)$$

最好的参考向量是最小化总重构误差的参考向量。 b_i^t 也依赖 m_i ，并且我们不能解析地求解这个优化问题。对此，我们有一个称作 k -均值聚类 (k -means clustering) 的迭代过程：首先，我们以随机初始化的 m_i 开始。然后，在每次迭代中，我们先对每个 x^t ，使用式 (7.4) 计算估计标号 (estimated labels) b_i^t ；如果 b_i^t 为 1，则 x^t 属于分组 m_i 。然后，一旦我们有了这些标号，我们就最小化式 (7.3)。取它关于 m_i 的导数并令其等于 0，我们得到

$$m_i = \frac{\sum_t b_i^t x^t}{\sum_t b_i^t} \quad (7.5)$$

参考向量被设置为它所代表的所有实例的均值。注意，除了用估计的标号 b_i^t 取代标号 r_i^t 外，这与式 (7.2) 中的均值公式相同。这是一个迭代过程，因为一旦我们计算了新的 m_i ， b_i^t

改变并且需要重新计算, 这反过来又影响 m_i 。这个两步过程一直重复, 直到 m_i 稳定(见图 7-2)。 k -均值算法的伪代码在图 7-3 中给出。

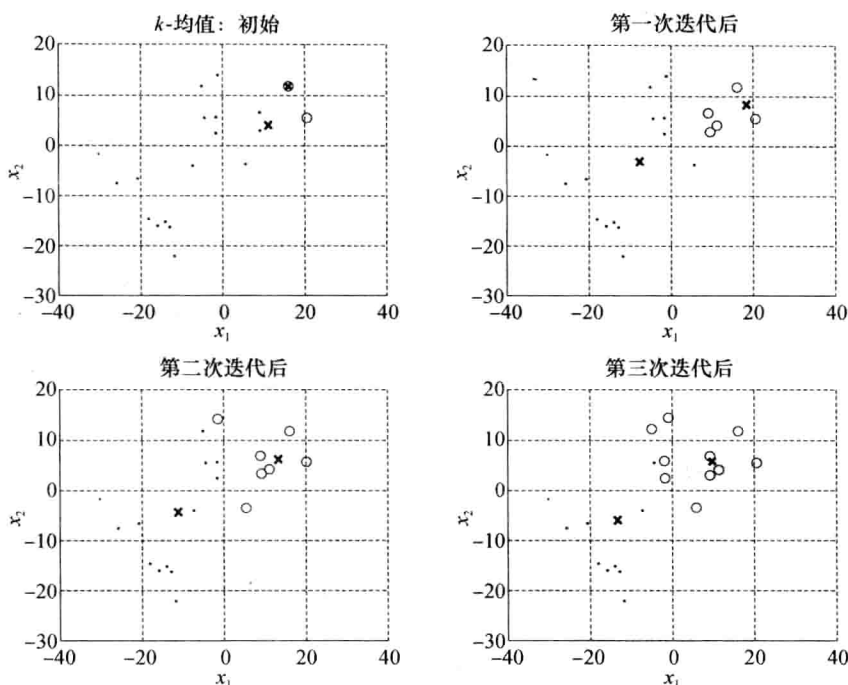


图 7-2 k -均值演变。叉指示中心位置。数据点根据最近的中心标记

```

初始化  $m_i, i = 1, \dots, k$ ; 例如, 将  $m_i$  初始化为  $k$  个随机的  $x^i$ 
Repeat
  For 所有的  $x^i \in X$ 
     $b_i^t \leftarrow \begin{cases} 1 & \text{如果 } \|x^i - m_i\| = \min_j \|x^i - m_j\| \\ 0 & \text{否则} \end{cases}$ 
  For 所有的  $m_i, i = 1, \dots, k$ 
     $m_i \leftarrow \sum_i b_i^t x^i / \sum_i b_i^t$ 
Until  $m_i$  收敛
  
```

图 7-3 k -均值算法

k -均值的一个缺点是它是一个局部搜索过程, 并且最终的 m_i 高度依赖于初始的 m_i 。对于初始化, 存在各种不同的方法:

- 可以简单地随机选择 k 个实例作为初始的 m_i 。
- 可以计算所有数据的均值, 并将一些小随机向量加到均值上, 得到 k 个初始的 m_i 。
- 可以计算主成分, 将它的值域化分成 k 个相等的区间, 将数据化分成 k 个分组, 然后取这些分组的均值作为初始中心。

收敛后, 所有的簇中心应当涵盖数据实例的某个子集并且是有用的, 因此, 最好将中心

初始化在有数据的地方。

还有一些算法动态地添加新中心或删除空的中心。在领导者聚类算法 (leader cluster algorithm) 中, 一个远离 (用一个阈值定义) 已有中心的实例导致在该点创建一个新中心 (我们将在第 12 章讨论这样一种神经网络算法 ART)。或者, 一个涵盖了大量实例 ($\sum_i b_i^t / N > \theta$) 的中心可以分裂为两个 (通过把一个小随机向量加到两个拷贝中的一个上, 使得它们不同)。类似地, 涵盖太少实例的中心可以删除, 并从输入空间的某个其他部分重新开始。

k -均值算法用于聚类, 也就是说用于找出数据的分组, 其中分组用它们的中心 (分组的典型代表) 表示。向量量化是聚类的一种应用, 但是聚类也用于分类或回归阶段之前的预处理。给定 \mathbf{x}^t , 当我们计算 b_i^t 时, 我们做了从原来的空间到 k -维空间的映射, 即映射到 k -维超立方体的一个隅角上。然后, 可以在这个新空间学习回归或判别式函数。我们将在第 12 章讨论这样的方法。

7.4 期望最大化算法

在 k -均值中, 我们把聚类看作是寻找最小化总重构误差的码本向量问题。本节中, 我们的方法是概率的, 并且我们寻找最大化样本似然的支密度参数。使用式 (7.1) 的混合模型, 样本 $\mathcal{X} = \{\mathbf{x}^t\}_t$ 的对数似然为

$$\begin{aligned}\mathcal{L}(\Phi | \mathcal{X}) &= \log \prod_t p(\mathbf{x}^t | \Phi) \\ &= \sum_t \log \sum_{i=1}^k p(\mathbf{x}^t | \mathcal{G}_i) P(\mathcal{G}_i)\end{aligned}\quad (7.6)$$

其中 Φ 包含先验概率 $P(\mathcal{G}_i)$ 和支密度 $p(\mathbf{x}^t | \mathcal{G}_i)$ 的有效统计量。不幸的是, 我们不能解析地求解参数, 而需要借助于迭代优化。

期望最大化 (Expectation-Maximization, EM) 算法 (Dempster、Laird 和 Rubin 1977, Redner 和 Walker 1984) 用于最大似然估计, 其中问题涉及两组随机变量, 其中一组 X 是可观测的, 另一组 Z 是隐藏的。算法的目标是找到参数向量 Φ , 它最大化 X 的观测值的似然 $\mathcal{L}(\Phi | \mathcal{X})$ 。但是, 在不可行时, 我们关联附加的隐藏变量 (hidden variable) Z , 并使用二者表示潜在的模型, 最大化 X 和 Z 联合分布的似然, 完全 (complete) 似然 $\mathcal{L}_c(\Phi | \mathcal{X}, \mathcal{Z})$ 。

由于 Z 值不是观测的, 我们不能直接求解完全数据似然 \mathcal{L}_c 。而是给定 \mathcal{X} 和当前参数值 Φ^l (其中 l 是迭代次数), 我们求它的期望 \mathcal{Q} 。这是算法的期望 (E) 步。然后, 在最大化 (M) 步, 我们寻找新的参数值 Φ^{l+1} , 它最大化期望。这样

E 步: $\mathcal{Q}(\Phi | \Phi^l) = E[\mathcal{L}_c(\Phi | \mathcal{X}, \mathcal{Z}) | \mathcal{X}, \Phi^l]$

M 步: $\Phi^{l+1} = \arg \max_{\Phi} \mathcal{Q}(\Phi | \Phi^l)$

Dempster、Laird 和 Rubin (1977) 证明增加 \mathcal{Q} 意味着增加不完全似然

$$\mathcal{L}(\Phi^{l+1} | \mathcal{X}) \geq \mathcal{L}(\Phi^l | \mathcal{X})$$

在混合模型的情况下, 隐藏变量是观测的源, 即哪个观测属于哪个分支。如果这些被给定, 例如, 作为监督学习的类标号被给定, 我们就想知道调整哪些参数, 以便拟合数据点。

146
149

EM 方法的执行过程如下：在 E 步，给定分支的当前知识，我们估计这些标号；而在 M 步，给定 E 步估计的标号，我们更新我们的类知识。这两步与 k -均值的两步相同： b_i^t 的计算(E 步)和 m_i 的重新估计(M 步)。

我们定义一个指示变量(indicator variable)向量 $\mathbf{z}^t = \{z_1^t, \dots, z_k^t\}$ ，其中如果 \mathbf{x}^t 属于簇 G_i ，则 $z_i^t = 1$ ，否则 z_i^t 为 0。 \mathbf{z} 是多项式分布，以先验概率 π_i 取自 k 个类，是 $P(G_i)$ 的简写。150

$$P(\mathbf{z}^t) = \prod_{i=1}^k \pi_i^{z_i^t} \quad (7.7)$$

观测 \mathbf{x}^t 的似然等于它的概率，被产生它的分支指定：

$$p(\mathbf{x}^t | \mathbf{z}^t) = \prod_{i=1}^k p_i(\mathbf{x}^t)^{z_i^t} \quad (7.8)$$

$p_i(\mathbf{x}^t)$ 是 $p(\mathbf{x}^t | G_i)$ 的简写。联合密度为

$$p(\mathbf{x}^t, \mathbf{z}^t) = P(\mathbf{z}^t) p(\mathbf{x}^t | \mathbf{z}^t)$$

而独立同分布的样本 χ 的完全数据似然为

$$\begin{aligned} \mathcal{L}_c(\Phi | \chi, \mathbf{Z}) &= \log \prod_t p(\mathbf{x}^t, \mathbf{z}^t | \Phi) \\ &= \sum_t \log p(\mathbf{x}^t, \mathbf{z}^t | \Phi) \\ &= \sum_t \log P(\mathbf{z}^t | \Phi) + \log p(\mathbf{x}^t | \mathbf{z}^t, \Phi) \\ &= \sum_t \sum_i z_i^t [\log \pi_i + \log p_i(\mathbf{x}^t | \Phi)] \end{aligned}$$

E 步：我们定义

$$\begin{aligned} \mathcal{Q}(\Phi | \Phi^t) &= E[\log P(\chi, \mathbf{Z}) | \chi, \Phi^t] \\ &= E[\mathcal{L}_c(\Phi | \chi, \mathbf{Z}) | \chi, \Phi^t] \\ &= \sum_t \sum_i E[z_i^t | \chi, \Phi^t] [\log \pi_i + \log p_i(\mathbf{x}^t | \Phi^t)] \end{aligned}$$

其中

$$\begin{aligned} E[z_i^t | \chi, \Phi^t] &= E[z_i^t | \mathbf{x}^t, \Phi^t] & \mathbf{x}^t \text{ 是独立同分布} \\ &= P(z_i^t = 1 | \mathbf{x}^t, \Phi^t) & z_i^t \text{ 是 0/1 随机变量} \\ &= \frac{p(\mathbf{x}^t | z_i^t = 1, \Phi^t) P(z_i^t = 1 | \Phi^t)}{p(\mathbf{x}^t | \Phi^t)} & \text{贝叶斯规则} \\ &= \frac{p_i(\mathbf{x}^t | \Phi^t) \pi_i}{\sum_j p_j(\mathbf{x}^t | \Phi^t) \pi_j} \\ &= \frac{p(\mathbf{x}^t | G_i, \Phi^t) P(G_i)}{\sum_j p(\mathbf{x}^t | G_j, \Phi^t) P(G_j)} \\ &= P(G_i | \mathbf{x}^t, \Phi^t) \equiv h_i^t \end{aligned} \quad (7.9)$$

我们看到隐藏变量的期望值 $E[z_i^t]$ 是 \mathbf{x}^t 被分支 G_i 产生的后验概率。因为这是概率，所以它在 0 和 1 之间，并且与 k -均值的 0/1“硬”标号不同，它是“软”标号。

M 步: 我们最大化 \mathcal{Q} , 得到下一组参数值 Φ^{l+1} :

$$\Phi^{l+1} = \arg \max_{\Phi} \mathcal{Q}(\Phi | \Phi^l)$$

其中

$$\begin{aligned} \mathcal{Q}(\Phi | \Phi^l) &= \sum_t \sum_i h_i^t [\log \pi_i + \log p_i(\mathbf{x}^t | \Phi^l)] \\ &= \sum_t \sum_i h_i^t \log \pi_i + \sum_t \sum_i h_i^t \log p_i(\mathbf{x}^t | \Phi^l) \end{aligned} \quad (7.10)$$

第二项独立于 π_i , 并且作为拉格朗日, 使用约束 $\sum_i \pi_i = 1$, 我们求解

$$\nabla_{\pi_i} \sum_t \sum_i h_i^t \log \pi_i - \lambda \left(\sum_i \pi_i - 1 \right) = 0$$

并且得到

$$\pi_i = \frac{\sum_t h_i^t}{N} \quad (7.11)$$

这类似于式(7.2)的先验计算。

类似地, 式(7.10)的第一项独立于分支, 并且可以在估计分支的参数时丢弃。我们解

$$\nabla_{\Phi} \sum_t \sum_i h_i^t \log p_i(\mathbf{x}^t | \Phi) = 0 \quad (7.12)$$

如果我们假定高斯分支 $\hat{p}_i(\mathbf{x}^t | \Phi) \sim \mathcal{N}(\mathbf{m}_i, \mathbf{S}_i)$, 则 M 步为

$$\begin{aligned} \mathbf{m}_i^{l+1} &= \frac{\sum_t h_i^t \mathbf{x}^t}{\sum_t h_i^t} \\ \mathbf{S}_i^{l+1} &= \frac{\sum_t h_i^t (\mathbf{x}^t - \mathbf{m}_i^{l+1})(\mathbf{x}^t - \mathbf{m}_i^{l+1})^T}{\sum_t h_i^t} \end{aligned} \quad (7.13)$$

152

这里, 对于 E 步的高斯分支, 我们计算

$$h_i^t = \frac{\pi_i |\mathbf{S}_i|^{-1/2} \exp[-(1/2)(\mathbf{x}^t - \mathbf{m}_i)^T \mathbf{S}_i^{-1}(\mathbf{x}^t - \mathbf{m}_i)]}{\sum_j \pi_j |\mathbf{S}_j|^{-1/2} \exp[-(1/2)(\mathbf{x}^t - \mathbf{m}_j)^T \mathbf{S}_j^{-1}(\mathbf{x}^t - \mathbf{m}_j)]} \quad (7.14)$$

又一次, 式(7.13)与式(7.2)之间的相似性并非偶然; 估计的软标号 h_i^t 取代了实际(未知的)标号 r_i^t 。

EM 用 k -均值初始化。在几次 k -均值迭代后, 我们得到中心 \mathbf{m}_i 的估计, 并且使用被每个中心涵盖的实例, 我们估计 \mathbf{S}_i 和 $\sum_i h_i^t/N$ 得到 π_i 。从那之后, 我们运行 EM, 如图 7-4 所示。

正如参数分类(见 5.5 节), 使用小样本和高维度, 我们可以通过化简假设来正则化。当 $\hat{p}_i(\mathbf{x}^t | \Phi) \sim \mathcal{N}(\mathbf{m}_i, \mathbf{S})$ 时, 在共享协方差矩阵的情况下, 式(7.12)化简为

$$\min_{\mathbf{m}_i, \mathbf{S}} \sum_t \sum_i h_i^t (\mathbf{x}^t - \mathbf{m}_i)^T \mathbf{S}^{-1} (\mathbf{x}^t - \mathbf{m}_i) \quad (7.15)$$

当 $\hat{p}_i(\mathbf{x}^t | \Phi) \sim \mathcal{N}(\mathbf{m}_i, s^2 \mathbf{I})$ 时, 在共享对角矩阵的情况下, 我们有

$$\min_{\mathbf{m}_i, s} \sum_t \sum_i h_i^t \frac{\|\mathbf{x}^t - \mathbf{m}_i\|^2}{s^2} \quad (7.16)$$

这是我们在 k -均值聚类中定义的重构误差式(7.3)。现在, 不同的是

$$h_i^t = \frac{\exp[-(1/2s^2) \|\mathbf{x}^t - \mathbf{m}_i\|^2]}{\sum_j \exp[-(1/2s^2) \|\mathbf{x}^t - \mathbf{m}_j\|^2]} \quad (7.17)$$

是 0 和 1 之间的概率。 k -均值聚类中的 b_i^j 做 0/1 硬决策, 而 h_i^j 是软标号, 它以一定概率将输入指派到簇中。当使用 h_i^j 而不是 b_i^j 时, 实例对所有分支的参数更新都有贡献, 对每个分支以一定的概率。当实例靠近两个中心的中点时, 这特别有用。这样, 我们看到 k -均值聚类是 EM 用于高斯混合模型的特例, 假定输入是独立的、具有相等和共享的方差, 并且标号是“硬的”。 k -均值用圆覆盖输入密度, 而 EM 一般用任意形状和任意方向的椭圆。

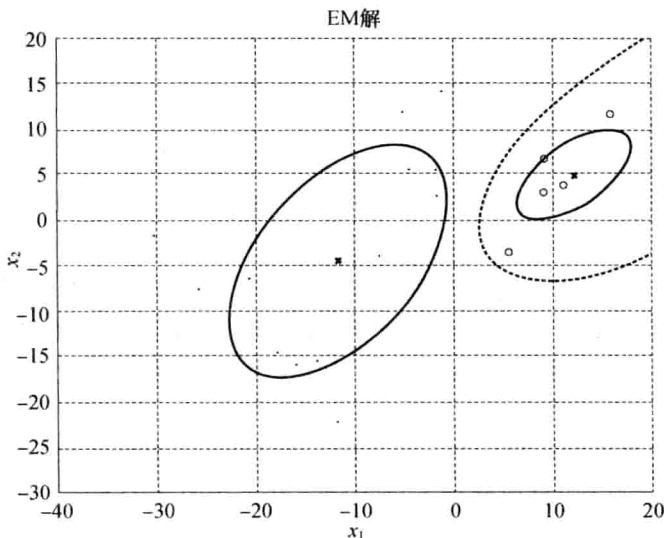


图 7-4 数据点和被 EM 拟合的高斯分布, 被图 7-2 的一个 k -均值迭代初始化。不像 k -均值, EM 允许估计协方差矩阵。图中显示了被较大的 h_i 标记的数据点、估计的高斯密度的围线和 $h_i = 0.5$ 的分离曲线(虚线)

7.5 潜在变量混合模型

当全协方差矩阵与高斯混合分布一起使用时, 即使没有奇异性, 如果输入维度很高而样本很小, 仍然有过拟合的危险。为了减少参数的个数而假定具有共同的协方差矩阵可能并不正确, 因为簇实际上可能具有不同的形状。假定对角矩阵可能更危险, 因为这排出了所有的相关性。另一种选择是在簇中做维归约。这减少了参数个数, 但仍然捕获相关性。自由参数的数量通过约化空间的维度控制。

当我们在簇中做因子分析时, 我们寻找产生簇中数据的潜在变量(latent variable)或隐藏变量(hidden variable)或因子(factor) (Bishop 1999):

$$p(\mathbf{x}^i | \mathcal{G}_i) \sim \mathcal{N}(\mathbf{m}_i, \mathbf{V}_i \mathbf{V}_i^T + \boldsymbol{\Psi}_i) \quad (7.18)$$

其中 \mathbf{V}_i 和 $\boldsymbol{\Psi}_i$ 是簇 \mathcal{G}_i 的因子载荷和特定方差。Rubin 和 Thayer (1982) 给出了因子分析的 EM 方法。可以把它推广到混合模型, 找到混合因子分析方法 (mixtures of factor analyzers) (Ghahramani 和 Hinton 1997)。在 E 步, 在式(7.9)中我们使用式(7.18); 而在 M 步, 我们对 \mathbf{V}_i 和 $\boldsymbol{\Psi}_i$ 而不是对 \mathbf{S}_i 求解方程(7.12)。类似地, 也可以在分组中做 PCA, 称作混合概率主成分分析方法 (mixtures of probabilistic principal component analyzer) (Tipping 和 Bishop

1999)。

当然,我们可以使用 EM 学习 S_i , 然后分别在每个簇上做 FA 或 PCA, 但是做 EM 更好, 因为它将两个步骤结合在一起, 并做软划分。每个实例对所有分组的潜在变量的计算都有贡献, 权重为 h'_i 。

7.6 聚类后的监督学习

像第 6 章讨论的维度归约方法一样, 聚类可以用于两个目的: 它可以用来探查数据, 理解数据的结构。维度归约方法用来发现变量之间的相关性, 从而对变量分组。聚类方法用来发现实例之间的相似性, 从而对实例分组。

如果找到这样的分组, 可以(通过领域专家)对它们命名, 定义它们的属性。我们可以选组均值作为分组中实例的代表原型, 或者可以写出属性的可能值域。这可以更简单地描述数据。例如, 如果一个公司的顾客看上去都归属于 k 个分组之一(顾客通过人口统计属性和与公司的交易勾画), 则将更好地理解顾客要素, 使得公司可以对不同类型的顾客使用不同的营销策略, 这是客户关系管理(customer relationship management, CRM)的一部分。同样, 公司还可以为没有归于大分组的顾客, 或需要特殊关注的顾客(例如, 定期买卖的顾客)制定策略。

聚类也常常作为预处理步骤使用。正如第 6 章的维度归约使得我们可以映射到新空间一样, 聚类后, 我们也可以映射到新的 k 维空间, 其中维是 h_i (或 b_i , 如果损失一些信息的话)。在监督学习的背景下, 我们可以在新空间学习判别式或回归函数。然而, 与诸如 PCA 等维度归约方法不同, 新空间的维度 k 可能比原空间的维度 d 大。

当我们使用像 PCA 这样的方法时, 新的维是原来的维的组合, 在新空间表示任意实例, 所有的维都有贡献, 即所有的 z_j 非零。在使用像聚类这种方法的情况下, 新的维是局部定义的, 存在很多新维 b_j , 但是它们之中只有一个(或几个, 如果使用 h_j)具有非零值。在前一种情况下, 存在少量维, 但都对表示起作用; 在后一种情况下, 存在许多维, 但只有少量起作用。

在监督学习前面使用非监督聚类或维度归约的优点是, 后者不需要标记的数据。标记数据的开销很大。我们可以使用大量未标记的数据学习簇参数, 然后使用少量标记的数据在第二阶段学习分类或回归。非监督学习又称“学习通常发生的事”(Barrow 1989)。当后随监督学习时, 我们先学习通常发生的事, 然后学习它意味什么。我们将在第 12 章讨论这种方法。

对于分类, 当每个类都是一个由大量分支组成的混合模型时, 整个密度是混合的混合密度(mixture of mixtures):

$$p(\mathbf{x} | C_i) = \sum_{j=1}^{k_i} p(\mathbf{x} | \mathcal{G}_{ij}) P(\mathcal{G}_{ij})$$

$$p(\mathbf{x}) = \sum_{i=1}^K p(\mathbf{x} | C_i) P(C_i)$$

其中 k_i 是组成 $p(\mathbf{x} | C_i)$ 的分支数, \mathcal{G}_{ij} 是类 i 的分支 j 。正如我们前面所讨论的, 分别为每个类学习分支的参数(或许在正则化之后)。这比用许多分支拟合所有的类的数据, 然后用类标记它们的方法好。

7.7 层次聚类

我们从统计学观点讨论了聚类，将聚类看作用一个混合模型拟合数据，或找出最小化重构误差的码字。还有一些聚类方法，它们只使用实例之间的相似性，而对数据没有其他要求；目标是找出分组，使得在同一个分组中的对象比在不同分组中的对象更相似。这种方法通过层次聚类(hierarchical clustering)实现。

这需要使用定义在实例间的相似性度量，或等价地，定义距离度量。通常使用欧氏距离，其中我们需要确保所有的属性都具有相同的尺度。欧氏距离是闵可夫斯基距离(Minkowski distance)的特例，其中 $p = 2$ ：

$$d_m(\mathbf{x}^r, \mathbf{x}^s) = \left[\sum_{j=1}^d (x_j^r - x_j^s)^p \right]^{1/p}$$

城市块距离(city-block distance)容易计算

$$d_{cb}(\mathbf{x}^r, \mathbf{x}^s) = \sum_{j=1}^d |x_j^r - x_j^s|$$

凝聚聚类(agglomerative clustering)算法从 N 个分组开始，每个分组最初只包含一个训练实例，重复合并相似的分组形成较大的分组，直到只有一个分组。分裂聚类(divisive clustering)以相反的方向，从单个分组开始，并将较大的分组合裂成较小的分组，直到每个分组包含单个实例。

在凝聚算法的每次迭代，我们选择两个最近的分组合并。在单链接聚类(single-link clustering)中，距离定义为两个分组的所有可能元素对之间的最小距离：

$$d(\mathcal{G}_i, \mathcal{G}_j) = \min_{\mathbf{x}^r \in \mathcal{G}_i, \mathbf{x}^s \in \mathcal{G}_j} d(\mathbf{x}^r, \mathbf{x}^s) \quad (7.19)$$

考虑一个加权的完全连接的图，顶点对应实例，顶点之间的边的权重对应实例之间的距离。单链接方法对应构造该图的最小生成树。

在全链接聚类(complete-link clustering)中，两个分组之间的距离取所有可能对之间的最大距离：

$$d(\mathcal{G}_i, \mathcal{G}_j) = \max_{\mathbf{x}^r \in \mathcal{G}_i, \mathbf{x}^s \in \mathcal{G}_j} d(\mathbf{x}^r, \mathbf{x}^s) \quad (7.20)$$

这两种是最频繁使用的、用于选择最近的分组合并的度量。其他可能的选择是使用所有可能点对之间平均距离的平均链接方法，度量两个分组成心(均值)之间距离的成心距离。

一旦运行了凝聚方法，结果通常被绘制成一个称作系统树图(dendrogram)的层次结构。这是一棵树，其中树叶对应实例，按照它们合并的次序分组。一个例子在图 7-5 中给出。树可以在任意水平截断，得到期望个数的分组。

单链接和全链接方法用不同的方法计算分组之间的距离，这影响聚类结果和系统树图：在单链接方法中，两个实例在水平 h 合并为一组，如果它们之间的距离小于 h ，或者存在一个中间实例序列使得相继实例之间的距离小于 h 。另一方面，在全链接中，一个分组中的所有实例之间的距离都小于 h 。单链接簇可能因这种“链接”效应而拉长。(在图 7-5 中，如果在 e 和 c 中间有一个实例会怎么样?)全链接簇趋向于更紧凑。

156
157

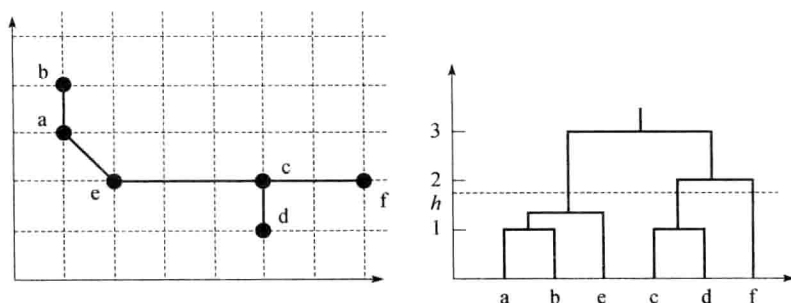


图 7-5 二维数据集和展示单链接聚类结果的系统树图。注意，树的树叶被排序使得分支不交叉。树在期望的值 h 上截断以得到簇

7.8 选择簇个数

像其他学习方法一样，聚类也有自己的调整复杂度的控制参数；这就是簇数 k 。给定 k ，聚类总是找出 k 个中心，不管它们是实际上有意义的分组，还是我们使用的算法强加的分组。存在多种调整 k 的方法：

- 在某些诸如颜色量化的应用中， k 由应用确定。
- 使用 PCA 在二维平面绘制数据可能用来发现数据的结构和数据中的簇数。
- 增量方法可能有助于确定 k ：设置允许的最大距离等价于设置每个实例的允许的最大重构误差。
- 在某些实际应用中，分组确认可以人工地进行，即检查簇是否实际上对数据中有意义的分组编码。例如，在数据挖掘应用中，领域专家可以做这项工作。在颜色量化中，我们可以目视检查图像，检查它的质量（尽管我们的眼睛和大脑并不逐个像素分析图像）。

依赖于我们使用的聚类方法类型，我们可以将重构误差或对数似然作为 k 的函数绘制图形，并找出“拐点”。足够大的 k 之后，算法将开始分裂分组，在这种情况下，重构误差将不会大幅度降低，对数似然将不会大幅度提高。类似地，在层次聚类，通过观察水平之间的差，我们可以决定好的划分。

7.9 注释

混合模型在统计学中被频繁使用。专门的教科书包括 Titterton、Smith 和 Makov (1985)，McLachlan 和 Basford (1988) 的书。McLachlan and Krishnan (1997) 讨论了 EM 算法的最近进展、如何加快它的收敛性和各种变形。在信号处理过程中， k -均值称作 *Linde-Buzo-Gray* (LBG) 算法 (Gersho 和 Gray 1992)。 k -均值频繁地用于统计学和信号处理的各种应用中，并且具有许多变形，其中之一是模糊 k -均值 (fuzzy k -means)。输入与分支的模糊关系也是一个 0 和 1 之间的数 (Bezdek 和 Pal 1995)。Alpaydm (1998) 比较了 k -均值、模糊 k -均值和高斯混合模型上的 EM。Xu 和 Jordan (1996) 给出了 EM 与学习高斯混合模型的其他学习算法的比较。在小样本上，另一种简化假设的方法是使用贝叶斯方法 (Ornoneit 和 Tresp 1996)。

Moerland(1999)在一组分类问题上比较了高斯混合模型和潜在变量混合模型,用实验展示了潜在变量模型的优点。Jain 和 Dubes(1988)是一本关于聚类的书,而 Jain、Murty 和 Flynn(1999)是一篇关于聚类的综述。

7.10 习题

1. 在图像压缩中, k -均值可以这样使用:图像被划分成 $c \times c$ 个窗口,并且这些 c^2 维向量构成样本。对于给定的 k (通常是 2 的幂),我们做 k -均值聚类。参考向量和每个窗口的下标通过通信线路发送。在接收端,通过使用下标读取参考向量表来重构图像。写一个计算机程序,对于不同的 k 和 c 值做这件事。对每种情况,计算重构误差和压缩率。
2. 我们可以做 k -均值聚类,划分实例,然后分别在每个分组计算 S_i 。为什么这不是一种好的想法?
3. 对于共享任意协方差矩阵 S 式(7.15)和 s^2 的情况,共享对角协方差矩阵式(7.16)的情况,为 S 推导 M 步公式。
4. 定义多元伯努利混合模型,其中输入是二元的,并推导 EM 公式。
5. 在分类混合的混合密度方法中,如何调整类 C_i 的分支数 k_i ?
6. 如何用二元输入向量进行层次聚类?例如,使用词袋表示对文本聚类。
7. 平均链接聚类与 k -均值聚类之间的相似和不同是什么?
8. 在层次聚类中,如何得到局部自适应距离?这样做的优缺点是什么?
9. 如何使得 k -均值对于离群点更鲁棒?
10. 产生系统树图后,能够对它“剪枝”吗?

160

7.11 参考文献

- Alpaydm, E. 1998. "Soft Vector Quantization and the EM Algorithm." *Neural Networks* 11: 467-477.
- Barrow, H. B. 1989. "Unsupervised Learning." *Neural Computation* 1: 295-311.
- Bezdek, J. C., and N. R. Pal. 1995. "Two Soft Relatives of Learning Vector Quantization." *Neural Networks* 8: 729-743.
- Bishop, C. M. 1999. "Latent Variable Models." In *Learning in Graphical Models*, ed. M. I. Jordan, 371-403. Cambridge, MA: MIT Press.
- Dempster, A. P., N. M. Laird, and D. B. Rubin. 1977. "Maximum Likelihood from Incomplete Data via the EM Algorithm." *Journal of Royal Statistical Society B* 39: 1-38.
- Gersho, A., and R. M. Gray. 1992. *Vector Quantization and Signal Compression*. Boston: Kluwer.
- Ghahramani, Z., and G. E. Hinton. 1997. *The EM Algorithm for Mixtures of Factor Analyzers*. Technical Report CRG TR-96-1, Department of Computer Science, University of Toronto (revised Feb. 1997).
- Jain, A. K., and R. C. Dubes. 1988. *Algorithms for Clustering Data*. New York: Prentice Hall.

- Jain, A. K., M. N. Murty, and P. J. Flynn. 1999. "Data Clustering: A Review." *ACM Computing Surveys* 31: 264-323.
- McLachlan, G. J., and K. E. Basford. 1988. *Mixture Models: Inference and Applications to Clustering*. New York: Marcel Dekker.
- McLachlan, G. J., and T. Krishnan. 1997. *The EM Algorithm and Extensions*. New York: Wiley.
- Moerland, P. 1999. "A Comparison of Mixture Models for Density Estimation." In *International Conference on Artificial Neural Networks*, ed. D. Willshaw, and A. Murray, 25-30. London, UK: IEE Press.
- Ornstein, D., and V. Tresp. 1996. "Improved Gaussian Mixture Density Estimates using Bayesian Penalty Terms and Network Averaging." In *Advances in Neural Information Processing Systems 8*, ed. D. S. Touretzky, M. C. Mozer, and M. E. Hasselmo, 542-548. Cambridge, MA: MIT Press.
- Redner, R. A., and H. F. Walker. 1984. "Mixture Densities, Maximum Likelihood and the EM Algorithm." *SIAM Review* 26: 195-239.
- Rubin, D. B., and D. T. Thayer. 1982. "EM Algorithms for ML Factor Analysis." *Psychometrika* 47: 69-76.
- Tipping, M. E., and C. M. Bishop. 1999. "Mixtures of Probabilistic Principal Component Analyzers." *Neural Computation* 11: 443-482.
- Titterton, D. M., A. F. M. Smith, and E. E. Makov. 1985. *Statistical Analysis of Finite Mixture Distributions*. New York: Wiley.
- Xu, L., and M. I. Jordan. 1996. "On Convergence Properties of the EM Algorithm for Gaussian Mixtures." *Neural Computation* 8: 129-151.
- Xu, R., and D. Wunsch II. 2005. "Survey of Clustering Algorithms." *IEEE Transactions on Neural Networks* 16: 645-678.

非参数方法

在前面的章节中，我们讨论了参数和半参数方法。在那里，我们假定数据取自一个形式已知的概率分布或混合分布。现在，我们将讨论非参数方法。当输入密度上不能做这样的假定时，可以使用非参数方法，并让数据自己说话。我们考虑密度估计、分类和回归的非参数方法，并讨论它们的时间和空间复杂度。

8.1 引言

在参数方法中，无论是密度估计、分类还是回归，我们都假设了一个在整个输入空间上有效的模型。例如，在回归中，当我们假定线性模型时，我们假定对于任何输入，输出都是输入的相同的线性函数。在分类中，当我们假定正态密度时，我们假定类的所有实例都取自这个相同的密度。参数方法的优点是，它把估计概率密度、判别式或回归函数问题归结为估计少量参数值。它的缺点是，假定并非总是成立的，并且不成立时可能导致很大的误差。

如果我们不能做这种假设并且不能使用参数模型时，一种可能的方法是使用如同我们在第7章所看到的半参数的混合模型，其中密度表示成几个参数模型的析取。在非参数估计(nonparametric estimation)中，我们只假定相似的输入具有相似的输出。这是一种合理的假设：世界是平稳的，并且无论是密度、判别式还是回归函数都缓慢地变化。相似的实例意味相似的事物。我们都爱我们的邻居，因为他们太像我们。

这样，我们的算法使用合适的距离度量，从训练集中找出相似的实例，并且由它们插值，得到正确的输出。不同的非参数方法在定义相似性或由相似的训练实例插值方法方面不同。在参数模型中，所有的训练实例都影响最终的全局估计。而在非参数情况下，不存在单个全局模型；需要时，局部模型被估计，它们只受邻近实例的影响。

在机器学习的文献中，非参数方法又称基于实例(instance-based)或基于记忆(memory-based)的学习算法，因为它们所做的是把训练实例存放在一个查找表中，并且由它们插值。这意味所有的训练实例都要存放，而存放所有训练实例需要 $\mathcal{O}(N)$ 存储量。此外，给定一个输入，应当找出相似的训练实例，而找出它们需要 $\mathcal{O}(N)$ 计算量。这种方法也称惰性(lazy)学习算法，因为不像急切(eager)的参数方法，当给定训练集时，它们并不计算模型，而是将模型的计算推迟到给定一个检验实例时才进行。对于参数学习方法，模型都相当简单，具有 $\mathcal{O}(d)$ 或 $\mathcal{O}(d^2)$ 量级个参数，并且一旦从训练集中计算出这些参数，我们保存模型并且在计算输出时就不再需要训练集了。通常， N 比 d (或 d^2)大得多，而这种存储和计算量的增加是非参数方法的缺点。

我们从估计密度函数开始，并讨论它在分类上的应用。然后，我们将该方法推广到回归中。

8.2 非参数密度估计

与通常的密度估计一样，我们假设样本 $\mathcal{X} = \{x^i\}_{i=1}^N$ 独立地从一个未知的概率密度 $p(\cdot)$ 中

抽取。 $\hat{p}(\cdot)$ 是 $p(\cdot)$ 的估计。我们从单变量情况开始,其中 x' 是标量,而稍后我们推广到多维情况。

累积分布函数 $F(x)$ 在点 x 的非参数估计是小于或等于 x 的样本所占的比例

$$\hat{F}(x) = \frac{\#\{x' \leq x\}}{N} \quad (8.1)$$

其中 $\#\{x' \leq x\}$ 表示其 x' 小于或等于 x 的训练样本数。类似地,密度函数的非参数估计可以用下式计算

$$\hat{p}(x) = \frac{1}{h} \left[\frac{\#\{x' \leq x+h\} - \#\{x' \leq x\}}{N} \right] \quad (8.2)$$

其中 h 是区间长度,并且假定落入该区间中的实例 x' 是“足够接近”的。本章提供的技术是一些变体,使用不同的启发式策略来确定邻近的实例和它们对估计的影响。

8.2.1 直方图估计

最古老、最流行的方法是直方图(histogram)。在直方图中,输入空间被划分成称作箱的相等区间。给定原点 x_0 和箱宽度 h ,箱是区间 $[x_0 + mh, x_0 + (m+1)h)$ (m 是正整数或负整数),而估计由下式给出

$$\hat{p}(x) = \frac{\#\{x' \text{ 与 } x \text{ 在相同的箱中}\}}{Nh} \quad (8.3)$$

在构造直方图时,我们必须选取原点和箱宽度。原点的选取影响靠近箱边界的估计,但影响估计的主要是箱宽度:使用小箱,估计是尖峰的,而使用大箱,估计较光滑(参见图8-1)。如果没有实例落入箱中,则估计为0,并且在箱边界处不连续。直方图的优点是:

一旦计算和存放了箱估计,我们就不再需要保留训练集。

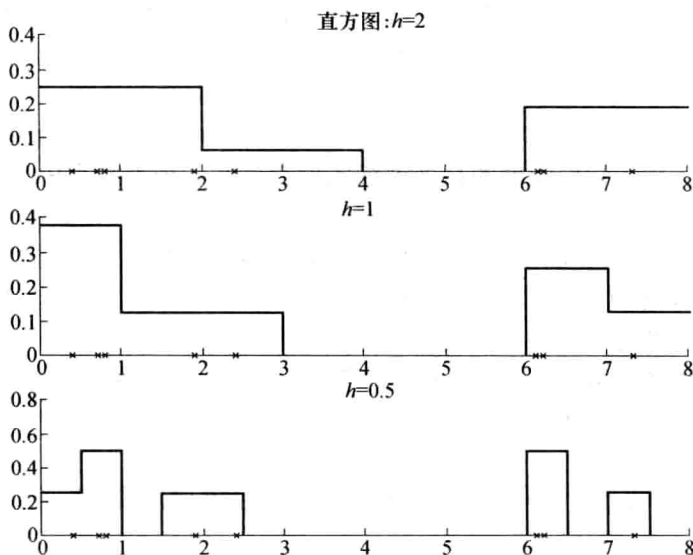


图8-1 各种箱长度的直方图。“×”表示数据点

质朴估计法 (naive estimator) (Silverman 1986) 使得我们不必设置原点。它定义为

$$\hat{p}(x) = \frac{\#\{x - h/2 < x^i \leq x + h\}}{2Nh} \quad (8.4)$$

并且它等于 x 总是落在宽度为 h 的箱中心的直方图估计 (参见图 8-2)。该估计还可以表示为

$$\hat{p}(x) = \frac{1}{Nh} \sum_{i=1}^N w\left(\frac{x - x^i}{h}\right) \quad (8.5)$$

其中权重函数定义为

$$w(u) = \begin{cases} 1 & \text{如果 } |u| < 1/2 \\ 0 & \text{否则} \end{cases}$$

这就好像每个 x^i 都有一个围绕它的大小为 h 的、对称的影响区域, 并且对落入该区域的 x 都贡献 1。于是, 非参数估计恰为其区域包含 x 的 x^i 的影响之和。因为这种影响区域是“硬的”(0 或 1), 所以估计不是连续函数并在 $x^i \pm h/2$ 处有跳跃。

166

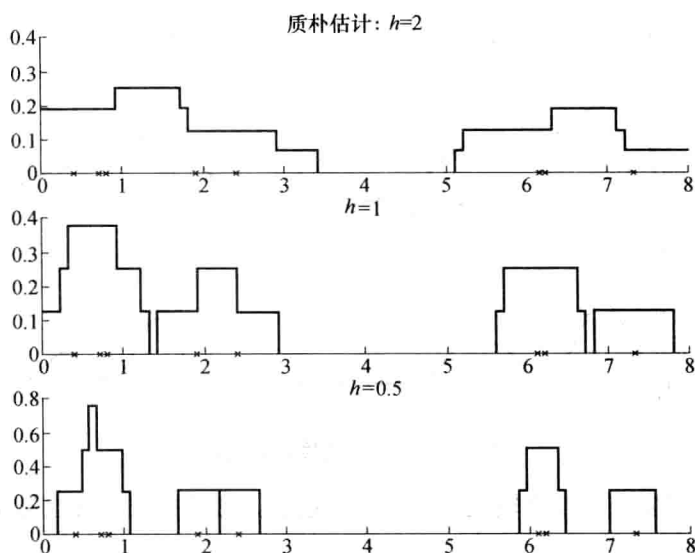


图 8-2 各种箱长度的质朴估计

8.2.2 核估计

为了得到光滑的估计, 我们使用一个光滑的权重函数, 称作核函数 (kernel function)。最流行的是高斯核:

$$K(u) = \frac{1}{\sqrt{2\pi}} \exp\left[-\frac{u^2}{2}\right] \quad (8.6)$$

核估计 (kernel estimator) 又称 Parzen 窗口 (Parzen windows), 定义为

$$\hat{p}(x) = \frac{1}{Nh} \sum_{i=1}^N K\left(\frac{x - x^i}{h}\right) \quad (8.7)$$

核函数 $K(\cdot)$ 决定影响的形状, 而窗口宽度 h 决定影响的宽度。就像质朴估计是“箱”的和一样, 核估计是“凸块”的和。所有的 x^i 都对 x 上的估计具有影响, 并且其影响随

$|x - x'|$ 的增加而平滑地减小。

167 为了简化计算, 如果 $|x - x'| > 3h$, 则 $K(\cdot)$ 可以取 0。还可以使用其他容易计算的核函数, 只要 $K(u)$ 对 $u=0$ 取最大值, 并且随 $|u|$ 增加而对称地减少。

当 h 很小时, 每个训练实例在一个小区域都具有较大影响, 而在较远的点上没有影响。当 h 较大时, 有更多的核重叠, 而我们得到较光滑的估计 (参见图 8-3)。如果 $K(\cdot)$ 处处非负并且积分为 1, 即如果它是合法的密度函数, 则 $\hat{p}(\cdot)$ 也是。此外, $\hat{p}(\cdot)$ 将继承核函数 $K(\cdot)$ 的连续性和可微性。例如, 如果 $K(\cdot)$ 是高斯函数, 则 $\hat{p}(\cdot)$ 将是光滑的并且具有所有导数。

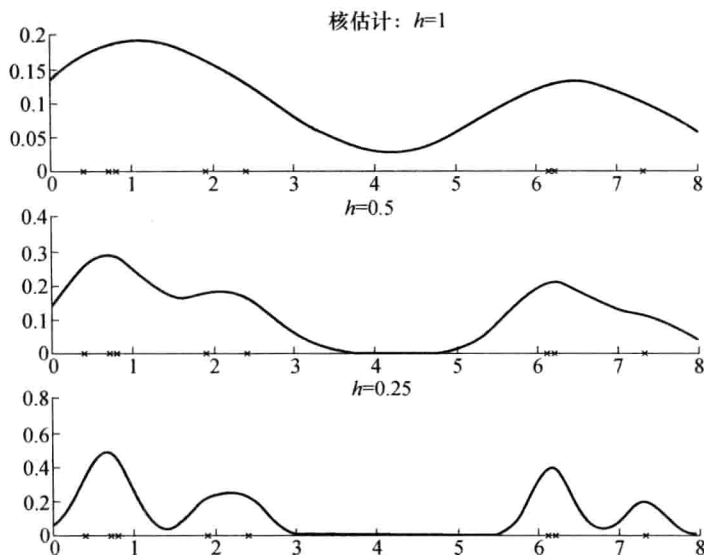


图 8-3 各种箱长度的核估计

一个问题是窗口宽度在整个输入空间上是固定的。已经提出各种自适应方法将 h 看作 x 周围密度的函数。

8.2.3 k -最近邻估计

168 估计的最近邻方法调整光滑量使之适应数据的局部密度。光滑度被所考虑的近邻数 k 控制。近邻数 k 远小于样本规模 N 。我们定义 a 和 b 之间的距离, 如为 $|a - b|$, 并对每个 x , 我们定义

$$d_1(x) \leq d_2(x) \leq \dots \leq d_N(x)$$

为从 x 到样本中的点按递增序排列的距离: $d_1(x)$ 是最远的样本的距离, $d_2(x)$ 是次近样本的距离, 如此下去。如果 x' 是数据点, 则我们定义 $d_1(x) = \min_i |x - x'|$, 并且如果 i 是最近样本的下标, 即 $i = \arg \min_i |x - x'|$, 则 $d_2(x) = \min_{j \neq i} |x - x'|$, 如此下去。

k -最近邻 (k -nearest neighbor, k -nn) 密度估计为

$$\hat{p}(x) = \frac{1}{2Nd_k(x)} \quad (8.8)$$

这就像 $h = 2d_k(x)$ 的质朴估计, 不同之处是我们不是固定 h 并检查多少样本落入箱中,

而是固定落入箱中的观测数 k ，并计算箱的大小。密度高的地方箱较小，而密度低的地方箱较大(参见图 8-4)。

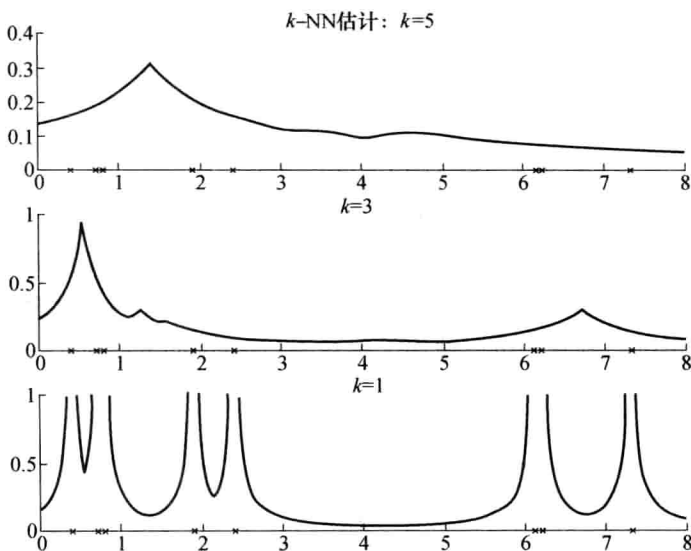


图 8-4 各种 k 值的 k -最近邻估计

k -nn 估计不是连续的；它的导数在所有的 $\frac{1}{2}(x^{(j)} + x^{(j+k)})$ 上不具有连续性，其中 $x^{(j)}$ 是样本的顺序统计量。 k -nn 不是概率密度函数，因为它的积分为 ∞ ，而不是 1。

为了得到更光滑的估计，我们可以使用其影响随距离增加而减小的核函数

$$\hat{p}(x) = \frac{1}{Nd_k(x)} \sum_{i=1}^N K\left(\frac{x - x^i}{d_k(x)}\right) \quad (8.9)$$

这就像具有自适应光滑参数 $h = d_k(x)$ 的核估计。通常， $K(\cdot)$ 取高斯核。

8.3 到多元数据的推广

给定 d -维观测的样本 $\mathcal{X} = \{\mathbf{x}^i\}_{i=1}^N$ ，多元核密度估计为

$$\hat{p}(\mathbf{x}) = \frac{1}{Nh^d} \sum_{i=1}^N K\left(\frac{\mathbf{x} - \mathbf{x}^i}{h}\right) \quad (8.10)$$

满足必要条件

$$\int_{\mathbb{R}^d} K(\mathbf{x}) d\mathbf{x} = 1$$

一个显然的候选是多元高斯核：

$$K(\mathbf{u}) = \left(\frac{1}{\sqrt{2\pi}}\right)^d \exp\left[-\frac{\|\mathbf{u}\|^2}{2}\right] \quad (8.11)$$

然而，由于维灾难 (curse of dimensionality)，在高维空间使用非参数估计时需要小心：令 \mathbf{x} 是 8 维的，我们使用每维 10 个箱的直方图，则有 10^8 个箱。除非我们有大量数据，否

则大部分箱为空, 并且那里的估计为 0。在高维空间, “近邻”概念也变得模糊不清, 因此在选择 h 时需要小心。

例如, 式(8.11)中欧几里得范数的使用意味核在所有维上都具有相等的尺度。如果输入具有不同的尺度, 则应当将它们规范化, 使其具有相同的方差。这还没有考虑相关性, 并且当核函数与潜在的分布具有相同形式时, 将获得更好的结果

$$K(\mathbf{u}) = \frac{1}{(2\pi)^{d/2} |\mathbf{S}|^{1/2}} \exp \left[-\frac{1}{2} \mathbf{u}^T \mathbf{S}^{-1} \mathbf{u} \right] \quad (8.12)$$

其中 \mathbf{S} 是样本协方差矩阵。这对应使用马氏距离而不是欧氏距离。

也可以使用局部距离度量, 其中 \mathbf{S} 由 \mathbf{x} 附近的实例计算, 例如由最近的 k 个实例计算。注意, 局部计算的 \mathbf{S} 可能是奇异的, 可能需要 PCA(或是 LDA, 在分类的情况中)。

如果输入是离散的, 则我们可以使用汉明距离 (Hamming distance), 它对不匹配的属性计数

$$HD(\mathbf{x}, \mathbf{x}') = \sum_{j=1}^d 1(x_j \neq x'_j) \quad (8.13)$$

其中

$$1(x_j \neq x'_j) = \begin{cases} 1 & \text{如果 } x_j \neq x'_j \\ 0 & \text{否则} \end{cases}$$

然后, 在使用 $\|\mathbf{x} - \mathbf{x}'\|$ 或 $(\mathbf{x} - \mathbf{x}')^T \mathbf{S}^{-1} (\mathbf{x} - \mathbf{x}')$ 的地方使用 $HD(\mathbf{x}, \mathbf{x}')$ 进行核估计或找出 k 个最近邻。

8.4 非参数分类

当用于分类时, 我们使用非参数方法估计类条件密度 $p(\mathbf{x} | C_i)$ 。类条件密度的核估计由下式给出

$$\hat{p}(\mathbf{x} | C_i) = \frac{1}{N_i h^d} \sum_{t=1}^N K\left(\frac{\mathbf{x} - \mathbf{x}^t}{h}\right) r_i^t \quad (8.14)$$

其中如果 $\mathbf{x}^t \in C_i$, 则 r_i^t 为 1, 否则 r_i^t 为 0。 N_i 是属于 C_i 的标记实例数; $N_i = \sum_t r_i^t$ 。先验密度的 MLE 是 $\hat{P}(C_i) = N_i/N$ 。于是, 判别式可以表示为

$$\begin{aligned} g_i(\mathbf{x}) &= \hat{p}(\mathbf{x} | C_i) \hat{P}(C_i) \\ &= \frac{1}{N h^d} \sum_{t=1}^N K\left(\frac{\mathbf{x} - \mathbf{x}^t}{h}\right) r_i^t \end{aligned} \quad (8.15)$$

并且 \mathbf{x} 被指派到判别式取最大值的类。公共因子 $1/(N h^d)$ 可以忽略。这样, 每个训练实例都为它的类投票, 而对其他类没有影响; 投票的权重由核函数 $K(\cdot)$ 给定, 通常赋予更近的实例更高的权重。

对于 k -nn 估计的特殊情形, 我们有

$$\hat{p}(\mathbf{x} | C_i) = \frac{k_i}{N_i V^k(\mathbf{x})} \quad (8.16)$$

其中 k_i 是 k 个最近邻中属于 C_i 的近邻数, 而 $V^k(\mathbf{x})$ 是中心在 \mathbf{x} , 半径为 $r = \|\mathbf{x} - \mathbf{x}_{(k)}\|$ 的

d -维超球的体积, 这里 $\mathbf{x}_{(k)}$ 是(所有类 \mathbf{x} 的近邻中)第 k 个距离 \mathbf{x} 最近的观测: $V^* = r^d c_d$, c_d 是 d 维单位球的体积。例如, $c_1 = 2$, $c_2 = \pi$, $c_3 = 4\pi/3$, 如此等等。于是

$$\hat{P}(C_i | \mathbf{x}) = \frac{\hat{p}(\mathbf{x} | C_i) \hat{P}(C_i)}{\hat{p}(\mathbf{x})} = \frac{k_i}{k} \quad (8.17)$$

k -nn 分类法(k -nn classifier)将输入指派到输入的 k 个最近邻中具有最多实例的类。所有的近邻都有相同的投票权, 并且选取 k 个近邻中具有最多投票者的类。平局随意打破或用加权投票。通常, k 取奇数, 以减少平局: 难以区分的情况一般出现在两个相邻的类之间。

再则, 欧氏距离的使用对应假定不相关的输入具有相同的方差。如果不是这种情况, 应当选择合适的度量。一个例子是判别式自适应最近邻(discriminant adaptive nearest neighbor) (Hastie 和 Tibshirani 1996), 那里局部地估计将类分离的最佳距离。

k -nn 的一种特殊情况是最近邻分类(nearest neighbor classifier), 其中 $k = 1$, 并且输入被指派到最近的模式所在的类。这将空间划分成 Voronoi 图[⊖] (Voronoi tessellation) 形式(参见图 8-5)。

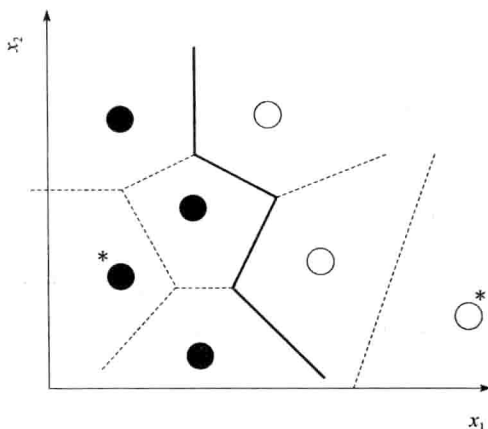


图 8-5 虚线是 Voronoi 图, 而实线是判别式。在精简的最近邻中, 可以删除那些不参与定义判别式的实例(用“*”标记)而不增加训练误差

8.5 精简的最近邻

非参数方法的时间和空间复杂度与训练集的大小成正比。已经提出了一些精简方法, 以减少存放的实例数而不降低性能。其基本思想是选择 \mathcal{X} 的最小子集 \mathcal{Z} 使得用 \mathcal{Z} 替代 \mathcal{X} 时, 误差不增加 (Dasarathy 1991)。

最著名和最早的方法是精简的最近邻(condensed nearest neighbor), 那里使用 1-nn 作为分类的非参数估计 (Hart 1968)。1-nn 以分段线形的方式逼近判别式, 并且只需要保存定义

⊖ Voronoi tessellation 又称 Voronoi diagram, 是 Georgy Voronoi 提出的。它由一组由连接两个相邻点线段的垂直平分线组成的多边形组成。——译者注

判别式的实例。类区域内部的实例不必作为它的同一类的最近邻存放，并且它的缺失不会导致(训练集上的)任何错误(见图8-5)。这样的子集称作相容子集，并且我们希望找出最小的相容子集。

Hart 提出了一种发现 Z 的贪心算法(见图8-6)：该算法从空集 Z 开始，以随机次序逐个扫描 X 中的实例，并检查它们是否能够被 1-nn 用已经在 Z 中的实例正确地分类。如果一个实例被错误分类，则将它添加到 Z 中；如果它被正确分类，则 Z 不变。应当扫描数据集多遍，直到没有实例再添加到 Z 中。该算法进行局部搜索，并且依赖于看到训练实例的次序，可能找出不同的子集，每个子集在验证数据上具有不同的准确率。因此，不能保证找到最小的相容子集。找出最小相容子集是 NP-完全问题(Wilfong 1992)。

```

 $Z \leftarrow \emptyset$ 
Repeat
  For 所有的  $x \in X$  (以随机次序)
    找出  $x' \in Z$  使得  $\|x - x'\| = \min_{x' \in Z} \|x - x'\|$ 
    If  $\text{class}(x) \neq \text{class}(x')$  将  $x$  添加到  $Z$  中
Until  $Z$  不改变
  
```

图 8-6 精简的最近邻算法

精简的最近邻是一种贪心的算法，旨在最小化训练误差和用存放的子集规模度量的复杂度。我们可以写一个增广的误差函数

$$E'(Z|X) = E(X|Z) + \lambda |Z| \quad (8.18)$$

其中 $E(X|Z)$ 是存放 Z 在 X 上的误差。 $|Z|$ 是 Z 的基数，而第二项是对复杂度加罚。与所有的正则化方案一样， λ 体现误差与复杂度之间的折中，使得对于较小的 λ ，误差变得更重要，并且随着 λ 增大，对复杂的模型加更大的罚。尽管精简的最近邻是一种最小化式(8.18)的方法，但是还可以设计优化它的其他算法。

8.6 非参数回归：光滑模型

在回归中，给定训练集 $X = \{x^i, r^i\}$ ，其中 $r^i \in \mathcal{R}$ ，我们假定

$$r^i = g(x^i) + \varepsilon$$

在参数回归，我们假定某次多项式，并计算它的系数，最小化训练集上误差的平方和。当不能假定这种多项式时，使用非参数回归；我们只假定相近的 x 具有相近的 $g(x)$ 值。与非参数密度估计一样，给定 x ，我们的方法是找到 x 的邻域，并求邻域中 r 的平均值，计算 $\hat{g}(x)$ 。非参数回归估计子又称光滑子(smooth)，而该估计称作光滑(Härdle 1990)。类似于密度估计，存在各种定义邻域和在邻域中取平均值的方法。我们对一元变量 x 讨论这些方法；与密度估计一样，使用多元核函数，可以用直截了当的方式把它们推广到多元情况中。

8.6.1 移动均值光滑

如果我们像在直方图中那样，定义一个原点和箱宽度并在箱中求 r 的平均值，则我们得到回归图(regressogram)(参见图8-7)。

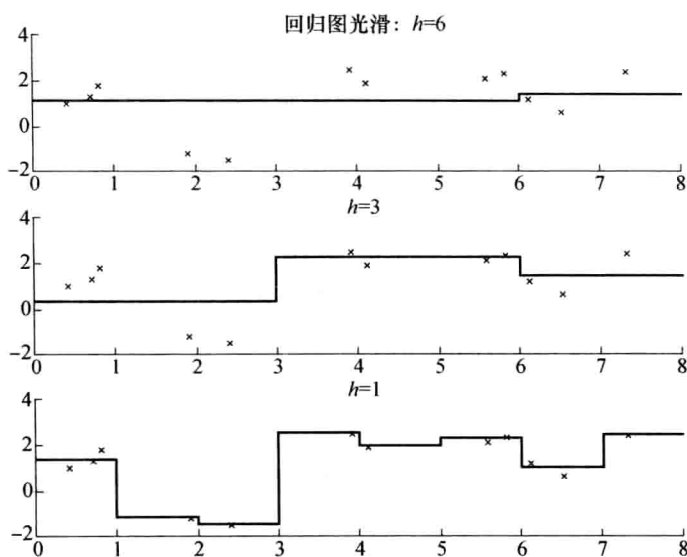


图 8-7 各种箱长度的回归图。“x”表示数据点

$$\hat{g}(x) = \frac{\sum_{i=1}^N b(x, x^i) r^i}{\sum_{i=1}^N b(x, x^i)} \quad (8.19)$$

其中

$$b(x, x^i) = \begin{cases} 1 & \text{如果 } x^i \text{ 与 } x \text{ 在同一个箱中} \\ 0 & \text{否则} \end{cases}$$

由于需要固定原点，箱边界上的不连续是令人烦恼的。像质朴估计一样，在移动均值光滑 (running mean smoother) 中，我们在 x 周围定义一个对称的箱并在那里取平均值 (见图 8-8)。

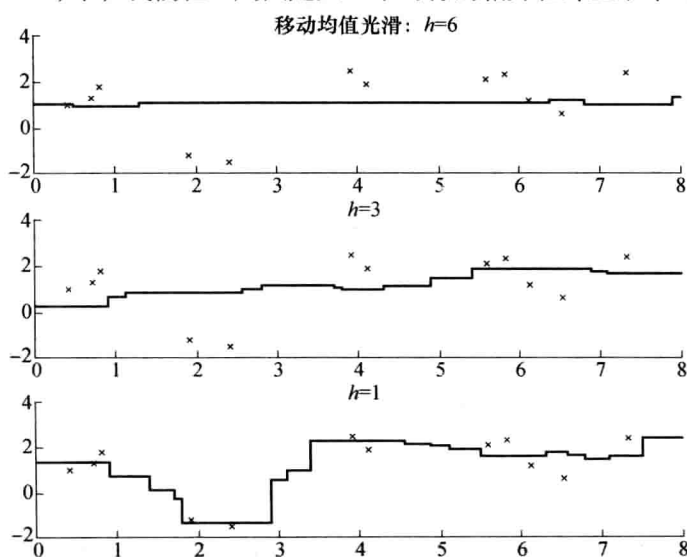


图 8-8 各种箱长度的移动均值光滑

$$\hat{g}(x) = \frac{\sum_{i=1}^N w\left(\frac{x-x^i}{h}\right) r^i}{\sum_{i=1}^N w\left(\frac{x-x^i}{h}\right)} \quad (8.20)$$

其中

$$w(u) = \begin{cases} 1 & \text{如果 } |u| < 1 \\ 0 & \text{否则} \end{cases}$$

这种方法在平滑分段数据(例如时间序列)方面特别流行。在有噪声的应用中,我们可以使用箱中 r^i 的中位数,而不是它们的均值。

8.6.2 核光滑

与核估计一样,我们可以使用赋予较远的点较小权重的核函数,并且得到核光滑(kernel smoother)(见图8-9):

$$\hat{g}(x) = \frac{\sum_i K\left(\frac{x-x^i}{h}\right) r^i}{\sum_i K\left(\frac{x-x^i}{h}\right)} \quad (8.21)$$

通常使用高斯核 $K(\cdot)$ 。替换固定 h , 我们可以固定近邻数 k , 使得估计自动适应 x 周围的密度, 并且得到 k -nn 光滑(k -nn smoother)。

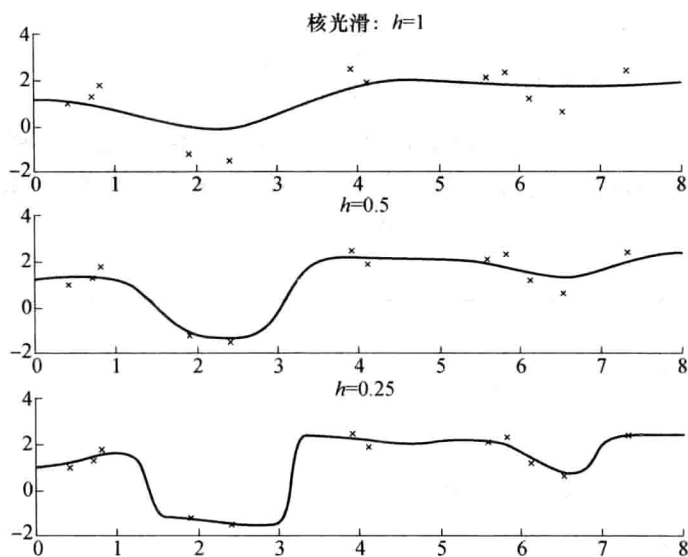


图8-9 各种箱长度的核光滑

8.6.3 移动线光滑

替代在点上取平均值和提供常量拟合, 我们可以对泰勒展开式多考虑一项并计算直线拟合。在移动线光滑(running line smoother)中, 我们可以使用邻域(被 h 或 k 定义)中的数据

点, 并拟合一个局部回归(见图 8-10)。

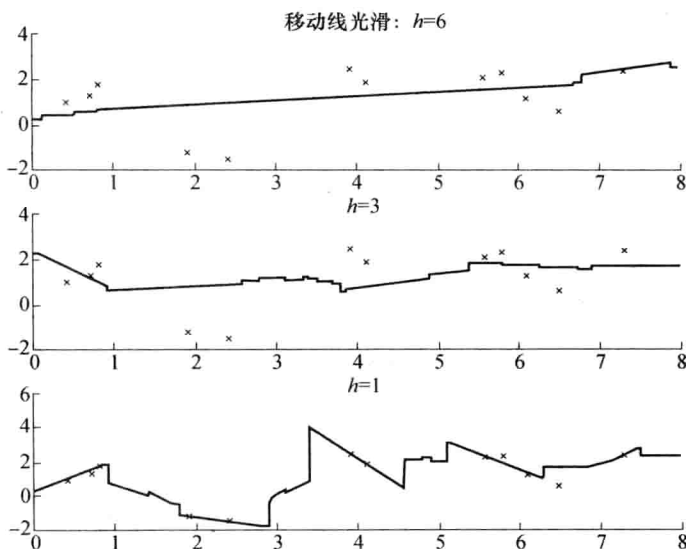


图 8-10 各种箱长度的移动线光滑

在局部加权移动线光滑(locally weighted running line smoother, 称作 loess[⊖])中, 我们使用核加权使得较远的点对误差具有较小影响, 而不是使用邻域的僵硬定义。

8.7 如何选择光滑参数

在非参数方法中, 对于密度估计或回归, 关键的参数是光滑参数, 如箱宽度或核扩展中的 h , 或近邻数 k 。目标是使得估计的不稳定性比数据点小。正如我们在前面已经讨论的, 数据中的易变性的一个根源是噪声, 其他根源是未知的潜在函数。我们应当光滑得恰好足以克服噪声——不少也不多。使用太大的 h 或 k , 许多实例都对点上的估计做出了贡献, 我们也光滑掉了源于函数的变化, 存在过分光滑。使用太小的 h 或 k , 单个实例具有很大影响, 我们甚至没有光滑掉噪声, 存在光滑不足。换句话说, 小 h 或 k 导致小偏倚但大方差。大 h 或 k 降低方差但增加偏倚。Geman、Bienenstock 和 Doursat(1992)讨论了非参数估计的偏倚和方差。

这一要求明确地表示在正如光滑样条(smoothing splines)中所使用的正则化函数中

$$\sum_i [r^i - \hat{g}(x^i)]^2 + \lambda \int_a^b [\hat{g}''(x)]^2 dx \quad (8.22)$$

第一项是拟合的误差。 $[a, b]$ 是输入区间; $\hat{g}''(\cdot)$ 是估计函数 $\hat{g}(\cdot)$ 的曲率(curvature), 因此度量变化。这样, 第二项处罚快速变化的估计。 λ 权衡变化和误差。例如, 使用大 λ , 我们得到更光滑的估计。

交叉验证用来调整 h , k 或 λ 。在密度估计中, 我们选择最大化验证集上似然的参数值。在监督环境下, 在训练集上试验一系列候选(见图 8-11), 选取最小化确认集上误差的参数值。

⊖ loess 意为局部回归。——译者注

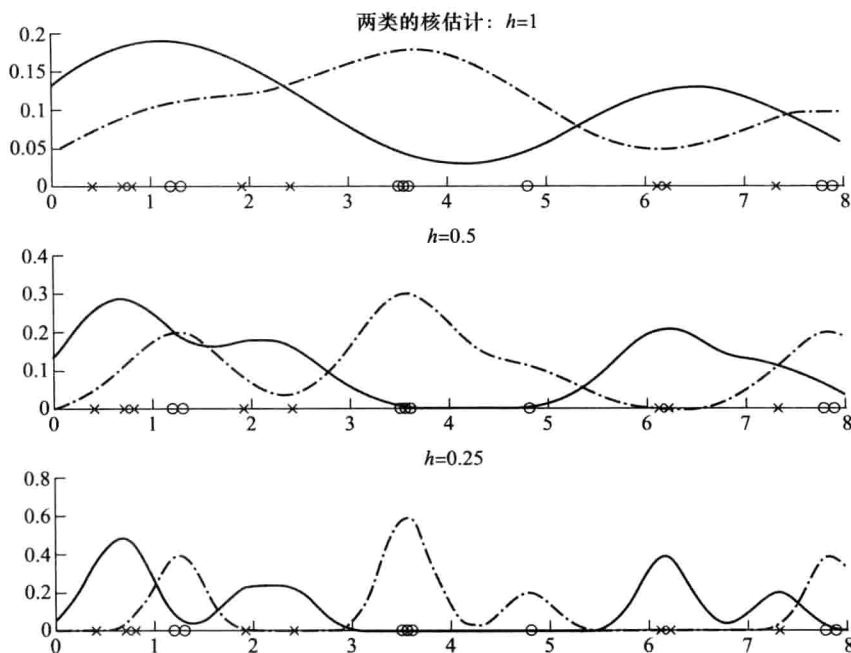


图 8-11 对于两类问题, 各种箱长度的核估计。绘制的是条件密度 $p(x | C_i)$, 顶部过光滑, 而底部欠光滑, 但是究竟哪个最好还依赖于验证数据指出

8.8 注释

k -最近邻和基于核的估计早在 50 年前就已提出了, 但是由于需要大量存储和计算, 该方法直到最近才流行 (Aha、Kibler 和 Albert 1991)。随着并行处理的进展, 存储和计算价格逐渐降低, 这些方法近来得到了更加广泛的使用。非参数估计方面的教科书是 Silverman 1986 和 Scott 1992。Dasarathy 1991 收集了许多关于 k -nn 和编辑/精简规则的文章。Aha 1997 是更加近期的工作的汇集。

非参数方法非常容易在单指令多数据 (SIMD) 机器上并行; 每个处理器在其局部存储器中存放一个实例, 而且并行地对该实例计算核函数值 (Stanfill 和 Waltz 1986)。乘以核函数可以看作卷积, 并且我们可以使用傅立叶变换更有效地计算估计 (Silverman 1986)。已经证明样条光滑等价于核光滑。

非参数估计最重要的因素是所使用的距离度量。对于离散属性, 我们可以简单地使用汉明距离, 那里我们只是累计非匹配的属性数。更复杂的距离函数在 Wettschereck, Aha 和 Mohri 1997 以及 Webb 1999 中被讨论。

在人工智能中, 非参数方法称作基于案例的推理 (case-based reasoning)。通过对已知的类似旧“案例”插值找到输出。这也允许知识提取: 给定的输出可以用列举这些类似的旧案例证明其合理性。

由于其简单性, k -nn 是最广泛使用的非参数分类方法, 并且在各种实践应用中相当成功。已经证明 (Cover 和 Hart 1967; 又见 Duda、Hart 和 Stork 2001): 在大样本中, 当 $N \rightarrow \infty$

时, 最近邻($k=1$)的风险不超过贝叶斯风险(我们能够得到的最好结果)的两倍, 并且从这方面来讲, 可以说“在被分类的无限样本集中, 一般的可用信息包含在最近邻中”(Cover 和 Hart 1967)。对于 k -nn, 业已证明, 随着 k 趋向于无穷大, 其风险逼近贝叶斯风险。

非参数回归在 Härdle 1990 中详细讨论。Hastie 和 Tibshirani (1990) 讨论了光滑模型并提出了加法模型(additive model), 其中多元函数被表示成一元估计的和。局部加权回归在 Atkeson, Moore 和 Schaal 1997 中讨论。这些模型与我们将在第 12 章讨论的径向基函数和混合专家模型很相似。

在精简的最近邻算法中, 我们看到只需要保存训练实例的一个子集, 这些实例靠近边界, 并且只使用它们就可以定义判别式。这一思想与我们将在第 13 章讨论的支持向量机(support vector machine)非常相似。第 13 章还将讨论度量实例之间各种核函数的相似性, 以及如何选择最好的核函数。把预测写成训练实例的影响和也构成了高斯过程(Gaussian process)的基础(参见第 14 章), 其中核函数称作协方差函数(covariance function)。

8.9 习题

1. 如何得到光滑的直方图?
2. 证明式(8.17)。
3. 如果 $k > 1$, 精简最近邻会怎么样?
4. 在精简的最近邻中, 先前添加到 Z 中的实例在之后的添加后可能不再是必需的。如何找出这种不再需要的实例?
5. 在回归图中, 替代箱中取平均值并做常量拟合, 我们可以使用落入箱中的实例并做线性拟合(参见图 8-12)。写出代码并与回归图做适当比较。

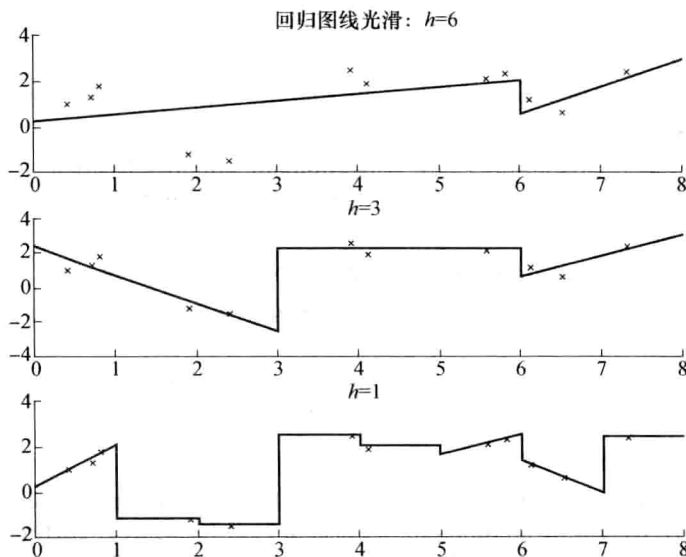


图 8-12 对于各种箱长度, 使用线性拟合的回归图

6. 为 8.6.3 节讨论的 loess 写出误差函数。

7. 提出一个移动均值估计的增量版本, 像精简最近邻一样, 它只在必要时存放实例。
8. 将核光滑推广到多元数据。
9. 在移动光滑中, 我们可以在检验点拟合一个常数、一条线或一个高阶多项式。如何在它们之间选择?
10. 在移动均值光滑中, 除了给出估计外, 还能计算出指示该点估计附近的方差(不确定性)的置信区间吗?

8.10 参考文献

- Aha, D. W., ed. 1997. Special Issue on Lazy Learning, *Artificial Intelligence Review* 11(1-5): 7-423.
- Aha, D. W., D. Kibler, and M. K. Albert. 1991. "Instance-Based Learning Algorithm." *Machine Learning* 6: 37-66.
- Atkeson, C. G., A. W. Moore, and S. Schaal. 1997. "Locally Weighted Learning." *Artificial Intelligence Review* 11: 11-73.
- Cover, T. M., and P. E. Hart. 1967. "Nearest Neighbor Pattern Classification." *IEEE Transactions on Information Theory* 13: 21-27.
- Dasarathy, B. V. 1991. *Nearest Neighbor Norms: NN Pattern Classification Techniques*. Los Alamitos, CA: IEEE Computer Society Press.
- Duda, R. O., P. E. Hart, and D. G. Stork. 2001. *Pattern Classification*, 2nd ed. New York: Wiley.
- Geman, S., E. Bienenstock, and R. Doursat. 1992. "Neural Networks and the Bias/Variance Dilemma." *Neural Computation* 4: 1-58.
- Härdle, W. 1990. *Applied Nonparametric Regression*. Cambridge, UK: Cambridge University Press.
- Hart, P. E. 1968. "The Condensed Nearest Neighbor Rule." *IEEE Transactions on Information Theory* 14: 515-516.
- Hastie, T. J., and R. J. Tibshirani. 1990. *Generalized Additive Models*. London: Chapman and Hall.
- Hastie, T. J., and R. J. Tibshirani. 1996. "Discriminant Adaptive Nearest Neighbor Classification." *IEEE Transactions on Pattern Analysis and Machine Intelligence* 18: 607-616.
- Scott, D. W. 1992. *Multivariate Density Estimation*. New York: Wiley.
- Silverman, B. W. 1986. *Density Estimation in Statistics and Data Analysis*. London: Chapman and Hall.
- Stanfill, C., and D. Waltz. 1986. "Toward Memory-Based Reasoning." *Communications of the ACM* 29: 1213-1228.
- Webb, A. 1999. *Statistical Pattern Recognition*. London: Arnold.
- Wettschereck, D., D. W. Aha, and T. Mohri. 1997. "A Review and Empirical Evaluation of Feature Weighting Methods for a Class of Lazy Learning Algorithms." *Artificial Intelligence Review* 11: 273-314.
- Wilfong, G. 1992. "Nearest Neighbor Problems." *International Journal on Computational Geometry and Applications* 2: 383-416.

决策树

决策树是一种实现分治策略的层次数据结构。它是一种有效的非参数学习方法，可以用于分类和回归。本章我们讨论由给定的标记的训练样本构造决策树的学习算法，以及如何将决策树转换成容易理解的简单规则的方法。

9.1 引言

对于参数估计，我们定义整个输入空间上的模型，并使用所有的训练数据学习它的参数。然后，对任意的检验输入，使用相同的模型和参数。对于非参数估计，我们把输入空间划分成被诸如欧几里得范数这样的距离度量定义的局部区域，并对每个输入使用由该区域的训练数据计算得到的对应的局部模型。在非参数模型中，给定一个输入，识别定义局部模型的局部数据的开销很大，需要计算从给定的输入到所有训练实例的距离。其计算复杂度为 $\mathcal{O}(N)$ 。

决策树(decision tree)是一种用于监督学习的层次模型，由此，局部区域通过少数几步递归分裂确定。决策树由一些内部决策节点和终端树叶组成(参见图9-1)。每个决策节点(decision node) m 实现一个具有离散输出的测试函数 $f_m(\mathbf{x})$ ，标记分支。给定一个输入，在每个节点应用一个测试，并根据测试的输出确定一个分支。这一过程从根节点开始，并递归地重复，直至到达一个树叶节点(leaf node)。这时，该树叶中的值形成输出。

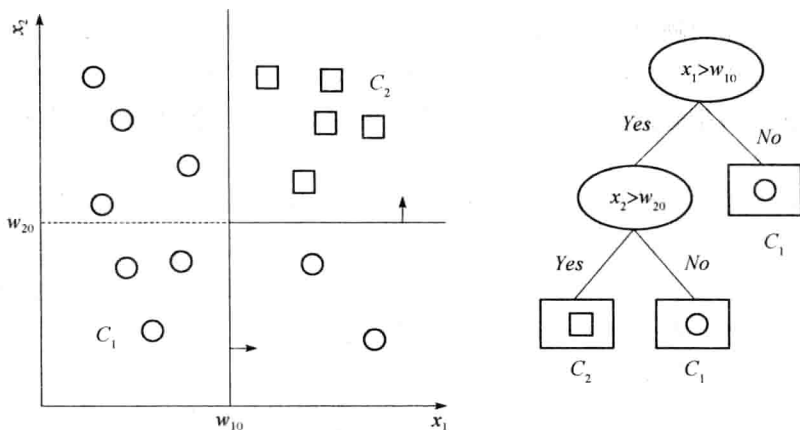


图9-1 数据集和对应的决策树。椭圆形节点是决策节点，而矩形节点是树叶节点。单变量的决策节点沿着一个轴划分，并且相继的划分相互正交。
第一次划分之后， $\{\mathbf{x} \mid x_1 < w_{10}\}$ 已是纯的，因此不需要再划分

每个 $f_m(\mathbf{x})$ 定义了一个 d -维输入空间中的判别式，将空间划分成较小区域。在从根节点

沿一条路径向下时, 这些较小的区域被进一步划分。 $f_m(\cdot)$ 是一个简单函数, 而作为树写下时, 复杂的函数被分解成一系列简单决策。不同的决策树方法对 $f_m(\cdot)$ 假设不同的模型, 而模型类确定了判别式的形状和区域的形状。每个树叶节点有一个输出标号。对于分类, 该标号是类代码; 而对于回归, 它是一个数值。一个树叶节点定义了输入空间的一个局部区域, 落入该区域的实例具有相同的输出。区域的边界被从树根到该树叶的路径上的内部节点中的判别式定义。

决策的层次安排使得涵盖输入的区域可以快速确定。例如, 如果决策是二元的, 则在最好情况下每个决策去掉一半实例。如果有 b 个区域, 则在最好情况下可以通过 $\log_2 b$ 次决策找到正确的区域。决策树的另一个优点是可解释性: 正如稍后我们将看到的, 可以把决策树转换成一组容易理解的 *IF-THEN* 规则。因此, 决策树非常流行, 并且常常比更准确但是不太好解释的方法更可取。

我们从一个决策节点只使用一个输入变量的单变量树开始, 考察如何为分类和回归构造这样的树。稍后, 我们将这种方法推广到一个内部节点可以使用所有输入的多变量树。

9.2 单变量树

在单变量树(univariate tree)中, 每个内部节点中的测试只使用一个输入维。如果所使用的输入维 x_j 是离散的, 取 n 个可能的值之一, 则该决策节点检查 x_j 的值, 并取相应的分支, 实现一个 n 路划分。例如, 如果属性是颜色, 具有可能的值{红, 蓝, 绿}, 则该属性上的节点具有三个分支, 每个对应该属性的三个可能值中的一个。

决策节点具有离散分支, 而数值输入应当离散化。如果 x_j 是数值的(有序的), 则测试是比较

$$f_m(\mathbf{x}) : x_j \geq w_{m0} \quad (9.1)$$

其中 w_{m0} 是适当选择的阈值。该决策节点将输入空间一分为二: $L_m = \{\mathbf{x} \mid x_j \geq w_{m0}\}$ 和 $R_m = \{\mathbf{x} \mid x_j < w_{m0}\}$; 称作一个二元划分(binary split)。从根到一个树叶的路径上的相继决策节点使用其他属性进一步把它们一分为二, 产生相互正交的划分。树叶节点定义输入空间中的超矩形(参见图9-1)。

树归纳是构造给定训练样本的树。对于给定的训练集, 存在许多对它进行无错编码的树, 而为了简单起见, 我们感兴趣的是寻找其中最小的树, 这里树的大小用树中的节点数和决策节点的复杂性度量。寻找最小树是 NP-完全的 (Quinlan 1986), 而我们必须使用基于启发式的局部搜索过程, 在合理的时间内得到合理的树。

树学习算法是贪心算法, 从包含全部训练数据的根开始, 每一步都选择最佳划分。依赖于所选取的属性是数值属性还是离散属性, 每次将数据划分成两个或 n 个子集。然后使用对应的子集递归地进行划分, 直到不再需要划分。此时, 创建一个树叶节点并标记它。

9.2.1 分类树

在用于分类的决策树, 即分类树(classification tree)中, 划分的优劣用不纯度度量(impurity measure)定量分析。一个划分是纯的, 如果对于所有分支, 划分后选择相同分支的所有

186
187

实例都属于相同的类。对于节点 m , 令 N_m 为到达节点 m 的训练实例数。对于根节点, N_m 为 N 。 N_m 个实例中 N_m^i 个属于 C_i 类, 而 $\sum_i N_m^i = N_m$ 。如果一个实例到达节点 m , 则它属于 C_i 类的概率估计为

$$\hat{p}(C_i | \mathbf{x}, m) \equiv p_m^i = \frac{N_m^i}{N_m} \quad (9.2)$$

节点 m 是纯的, 如果对于所有的 i , p_m^i 为 0 或 1。当到达节点 m 的所有实例都不属于 C_i 类时, p_m^i 为 0, 而当到达节点 m 的所有实例都属于 C_i 类时, p_m^i 为 1。如果划分是纯的, 则我们不需要进一步划分, 并可以添加一个树叶节点, 用 p_m^i 为 1 的类标记。一种度量不纯性的可能函数是熵函数(entropy) (Quinlan 1986) (参见图 9-2)。

$$I_m = - \sum_{i=1}^K p_m^i \log_2 p_m^i \quad (9.3)$$

其中 $0 \log 0 \equiv 0$ 。在信息论中, 熵是对一个实例的类代码进行编码所需要的最少位数。对于两类问题, 如果 $p^1 = 1$ 而 $p^2 = 0$, 则所有的实例都属于 C_1 类, 并且我们什么也不需要发送, 熵为 0。如果 $p^1 = p^2 = 0.5$, 则需要发送一位通告两种情况之一, 并且熵为 1。在这两个极端之间, 我们可以设计编码, 更可能的类用较短的编码, 更不可能的类用较长的编码, 每个信息使用不足一位。当存在 $K > 2$ 个类时, 相同的讨论成立, 并且当 $p^i = 1/K$ 时最大熵为 $\log_2 K$ 。

188

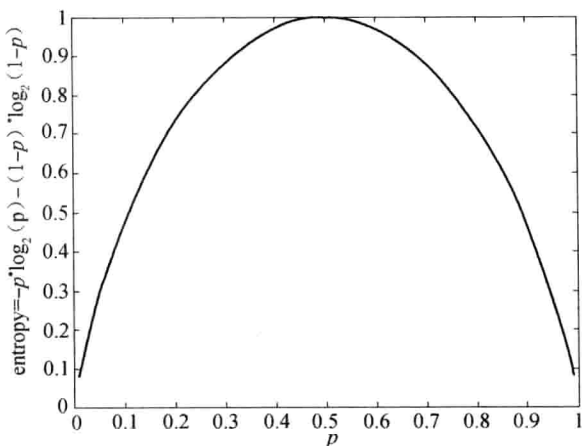


图 9-2 两类问题的熵函数

但是, 熵并非是唯一可能的度量。对于两类问题, 其中 $p^1 \equiv p$, $p^2 = 1 - p$, 函数 $\phi(p, 1 - p)$ 是非负函数, 度量划分的不纯度, 如果它满足如下性质 (Devroye、Györ 和 Lugosi 1996):

- 对于任意 $p \in [0, 1]$, $\phi(1/2, 1/2) \geq \phi(p, 1 - p)$ 。
- $\phi(0, 1) = \phi(1, 0) = 0$ 。
- 当 p 在 $[0, 1/2]$ 上时 $\phi(p, 1 - p)$ 是递增的, 而当 p 在 $[1/2, 1]$ 上时 $\phi(p, 1 - p)$ 是递减的。

例如

1) 熵

$$\phi(p, 1-p) = -p \log_2 p - (1-p) \log_2 (1-p) \quad (9.4)$$

式(9.3)是 $K > 2$ 个类的推广。

2) 基尼指数 (Gini index) (Breiman 等 1984)

$$\phi(p, 1-p) = 2p(1-p) \quad (9.5)$$

3) 误分类误差

$$\phi(p, 1-p) = 1 - \max(p, 1-p) \quad (9.6)$$

这些都可以推广到 $K > 2$ 类, 并且给定损失函数, 误分类误差可以推广到最小风险(习题1)。研究表明, 这三个度量之间并不存在显著差别。

如果节点 m 不是纯的, 则应当划分实例, 降低不纯度, 并且有多个属性可以用于划分。对于数值属性, 可能存在多个划分位置。在这些可能的划分中, 我们寻找最小化划分后的不纯度的划分, 因为我们希望产生最小的树。划分后的子集越纯, 则其后需要的划分(如果需要的话)就越少。当然, 这是局部最优, 并且不能保证找到最小的决策树。

设在节点 m , N_m 个实例中 N_{mj} 个取分支 j ; 这些是测试 $f_m(\mathbf{x}')$ 返回输出 j 的 \mathbf{x}' 。对于具有 n 个值的离散属性, 有 n 个输出; 而对于数值属性, 有两个输出 ($n=2$)。在两种情况下, 都满足 $\sum_{j=1}^n N_{mj} = N_m$ 。 N_{mj} 个实例中的 N_{mj}^i 个属于类 C_i : $\sum_{j=1}^K N_{mj}^i = N_{mj}$ 。类似地, $\sum_{j=1}^n N_{mj}^i = N_m^i$ 。

于是, 给定节点 m 、测试返回输出 j , 类 C_i 的概率估计为

$$\hat{p}(C_i | \mathbf{x}, m, j) \equiv p_{mj}^i = \frac{N_{mj}^i}{N_{mj}} \quad (9.7)$$

而划分后的总不纯度为

$$T_m = - \sum_{j=1}^n \frac{N_{mj}}{N_m} \sum_{i=1}^K p_{mj}^i \log p_{mj}^i \quad (9.8)$$

对于数值属性, 为了能够使用式(9.1)计算 p_{mj}^i , 我们还需要知道该节点的 w_{m0} 。在 N_m 个数据点之间, 存在 $N_m - 1$ 个可能的 w_{m0} : 我们不需要测试所有(无限多个)可能的点; 例如, 我们只需要考虑两点之间的中值就足够了。还要注意, 最佳划分总是在属于不同类的两个相邻点之间。这样, 我们检查每一个, 并取最高纯度作为该属性的纯度。对于离散属性, 不需要这种迭代。

对于所有的离散属性和数值属性, 对于数值属性的所有可能划分位置, 我们计算不纯度, 并且选取具有最小熵的划分位置, 例如在式(9.8)中度量的。于是, 对于所有的不纯的分支, 树构造递归地、平行地继续进行, 直到所有的分支都是纯的。这就是分类与回归树(classification and regression trees, CART)算法(Breiman 等, 1984)、ID3 算法(Quinlan 1986)和它的扩展 C4.5 (Quinlan 1993)的基本思想。算法的伪代码在图 9-3 中。

也可以说, 在树构造的每一步, 我们选择导致不纯度降低最多的划分。不纯度的降低是到达节点 m 的数据的不纯度(参见式(9.3))与划分后到达其分支的数据的总熵(参见式(9.8))之差。

189

190
191

```

Generate Tree( $\chi$ )
  If NodeEntropy( $\chi$ ) <  $\theta_t$  / * 式(9.3) */
    创建一个树叶, 用 $\chi$ 中的多数类标记
    Return
   $i \leftarrow \text{SplitAttribute}(\chi)$ 
  For  $x_i$  的每个分支
    找出落入该分支的 $\chi_i$ 
    Generate Tree( $\chi_i$ )

SplitAttribute( $\chi$ )
  MinEnt  $\leftarrow$  MAX
  For 所有的属性  $i = 1, \dots, d$ 
    If  $x_i$  是具有  $n$  个值的离散属性
      按照  $x_i$  将 $\chi$ 划分到 $\chi_1, \dots, \chi_n$ 
       $e \leftarrow \text{SplitEntropy}(\chi_1, \dots, \chi_n)$  / * 式(9.8) */
      If  $e < \text{MinEnt}$  MinEnt  $\leftarrow e$ ; bestf  $\leftarrow i$ 
    Else / *  $x_i$  是数值的 */
      For 所有可能的划分
        在  $x_i$  上将 $\chi$ 划分成 $\chi_1, \chi_2$ 
         $e \leftarrow \text{SplitEntropy}(\chi_1, \chi_2)$ 
        If  $e < \text{MinEnt}$  MinEnt  $\leftarrow e$ ; bestf  $\leftarrow i$ 
  Return bestf

```

图 9-3 构造分类树

一个问题是这种划分偏向于选择具有许多值的属性。当存在许多值时, 就存在许多分支, 并且不纯度可能很小。例如, 如果我们取训练样本的编号作为一个属性, 尽管它不是一个合理的特征, 但是不纯度度量将会选取它, 因为这样的话, 每个分支的不纯度都为 0。具有许多分支的节点是复杂的, 并且背离把类判别式划分成简单决策的思想。已经提出了许多方法对这样的属性加罚, 并权衡不纯度下降和分支因子两个因素。

179

当存在噪声时, 增长树直到最纯可能产生一棵非常大的、过分拟合的树。例如, 假设这种情况: 一个错误标记的实例混杂在一组正确标记的实例之中。为了减轻这种过分拟合, 当节点变得足够纯时, 树构造将终止; 即, 如果 $t < \theta_t$, 则数据子集就不再划分。这意味不需要使 p_{mj}^i 都恰为 0 或 1, 而只需要按照某个阈值 θ_p , p_{mj}^i 足够接近 0 或 1。在这种情况下, 创建一个树叶节点, 并将它标记为具有最大 p_{mj}^i 值的类。

θ_t (或 θ_p) 是复杂度参数, 与非参数估计中的 h 或 k 一样。当它们较小时, 方差大并且树增长较大, 以正确反映训练集; 而当它们较大时, 方差小并且树较小, 粗略地表示训练集并且可能具有较大偏倚。理想的值依赖于误分类的代价以及存储和计算开销。

一般地, 建议在树叶节点存放属于每个类的后验概率, 而不是用具有最大后验概率的类来标记树叶。这些概率在其后的步骤中可能是需要的。例如, 在计算风险时可能需要。注意, 我们不需要存放到达节点的实例或准确计数; 比率就足够了。

9.2.2 回归树

回归树(regression tree)可以用几乎与分类树完全相同的方法构造,唯一的不同是适合分类的不纯度度量用适合回归的不纯度度量取代。对于节点 m , 令 χ_m 为 χ 中到达节点 m 的子集, 即它是 $\mathbf{x} \in \chi$ 的满足从树根到节点 m 的所有决策节点条件的所有 \mathbf{x} 。我们定义

$$b_m(\mathbf{x}) = \begin{cases} 1 & \text{如果 } \mathbf{x} \in \chi_m: \mathbf{x} \text{ 到达结点 } m \\ 0 & \text{否则} \end{cases} \quad (9.9)$$

180 在回归树中, 划分的好坏用估计值的均方误差度量。令 g_m 为节点 m 中的估计值。

$$E_m = \frac{1}{N_m} \sum_t (r^t - g_m)^2 b_m(\mathbf{x}^t) \quad (9.10)$$

其中 $N_m = |\chi_m| = \sum_t b_m(\mathbf{x}^t)$ 。

在节点中, 我们使用到达该节点的实例要求的输出的均值(如果噪声太大用中值)

$$g_m = \frac{\sum_t b_m(\mathbf{x}^t) r^t}{\sum_t b_m(\mathbf{x}^t)} \quad (9.11)$$

于是, 式(9.10)对应于 m 上的方差。如果在一个节点上, 误差是可以接受的, 即 $E_m < \theta_r$, 则创建一个树叶节点, 存放 g_m 值。与第8章的回归图一样, 这会创建在叶边界不连续的分段常量近似。

如果误差不能接受, 则到达节点 m 的数据进一步划分, 使得诸分支的误差和最小。与分类一样, 在每个节点上, 我们寻找最小化误差的属性(和数值属性的划分阈值), 然后递归地进行上述过程。

令 χ_{mj} 为 χ_m 的取分支 j 的子集: $\bigcup_{j=1}^n \chi_{mj} = \chi_m$ 。我们定义

$$b_{mj}(\mathbf{x}) = \begin{cases} 1 & \text{如果 } \mathbf{x} \in \chi_{mj}: \mathbf{x} \text{ 到达结点 } m \text{ 并取分支 } j \\ 0 & \text{否则} \end{cases} \quad (9.12)$$

g_{mj} 是到达节点 m 的分支 j 的估计值。

$$g_{mj} = \frac{\sum_t b_{mj}(\mathbf{x}^t) r^t}{\sum_t b_{mj}(\mathbf{x}^t)} \quad (9.13)$$

而划分后的误差为

$$E'_m = \frac{1}{N_m} \sum_j \sum_t (r^t - g_{mj})^2 b_{mj}(\mathbf{x}^t) \quad (9.14)$$

192 对于任意划分, 误差的减少由式(9.10)和式(9.14)之差给出。我们寻找这样的划分, 它最大化误差的减少, 或等价地, 式(9.14)取最小值。将熵计算用均方误差替换, 类标号用平均值替换, 图9-3的程序代码可以用来训练回归树。

均误差是一种可能的误差函数; 另一种是最大误差

$$E_m = \max_j \max_t |r^t - g_{mj}| b_m(\mathbf{x}^t) \quad (9.15)$$

使用它，我们可以保证任意实例的误差都不大于给定的阈值。

可接受的误差阈值是复杂度的函数；其值越小，产生的树越大并且过分拟合的风险越大；其值越大，拟合不足和过分光滑的可能性越大(参见图 9-4 和图 9-5)。

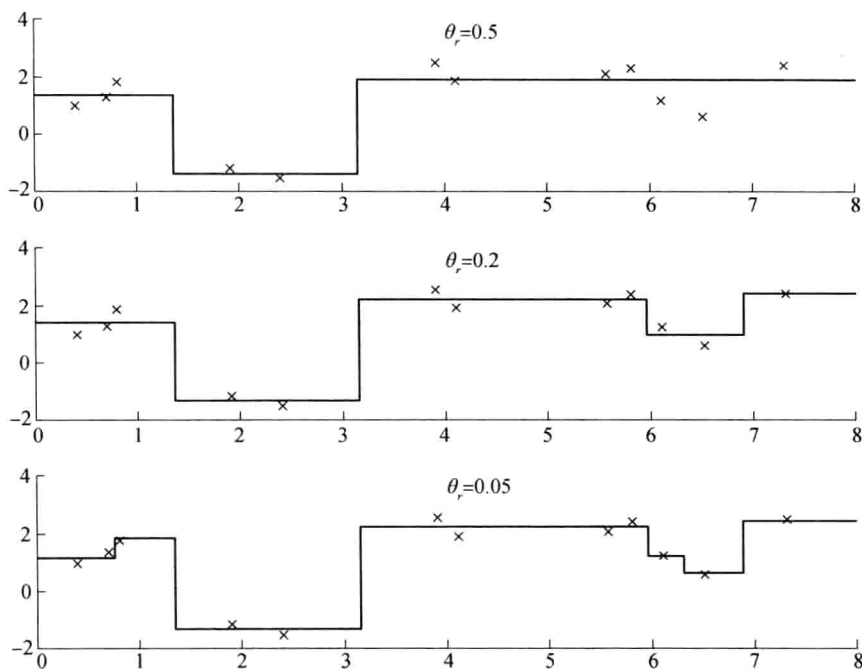


图 9-4 对于 θ_r 的不同值，回归树光滑。对应的树在图 9-5 中

类似于非参数回归中的从移动均值到移动直线，我们不是在实现常量拟合的树叶上取平均值，而是可以做线性回归拟合选定树叶上的实例：

$$g_m(\mathbf{x}) = \mathbf{w}_m^T \mathbf{x} + w_{m0} \quad (9.16)$$

这使得树叶上的估计依赖于 \mathbf{x} ，并且产生较小的树，但是这导致树叶节点上的附加的计算开销。

9.3 剪枝

通常，如果到达一个节点的训练实例数小于训练集的某个百分比(如 5%)，则无论是否不纯或是否有错误，该节点都不进一步分裂。其基本思想是：基于过少实例的决策树导致较大方差，从而导致较大泛化误差。在树完全构造出来之前提前停止树构造称作树的先剪枝(prepruning)。

得到较小树的另一种可能做法是后剪枝(postpruning)，实践中比先剪枝效果更好。前面，我们看到树的生长很贪心，在每一步，我们做出一个决策(即产生一个决策节点)并继续进行，绝不回溯尝试其他可能的选择。唯一的例外是后剪枝，那里我们试图找出并剪掉不必要的子树。

在后剪枝中，我们让树完全增长直到所有的树叶都是纯的并具有零训练误差。然后，我们找出导致过分拟合的子树并剪掉它们。我们从最初的被标记的数据集中保留一个剪枝集(pruning set)，在训练阶段不使用。对于每棵子树，我们用一个被该子树覆盖的训练实例标记的树叶节点替换它。如果该树叶在剪枝集上的性能不比该子树差，则剪掉该子树并保留树叶节点，因为子树的附加的复杂性是不必要的；否则保留子树。

例如, 在图 9-5 的第三棵树中, 有一个从条件 $x < 6.31$ 开始的子树。如果替换不增加剪枝集上的误差, 则该子树可以用树叶节点 $y = 0.9$ 替换(如第二棵树)。注意, 不要把剪枝集与确认集混淆, 它不同于确认集。

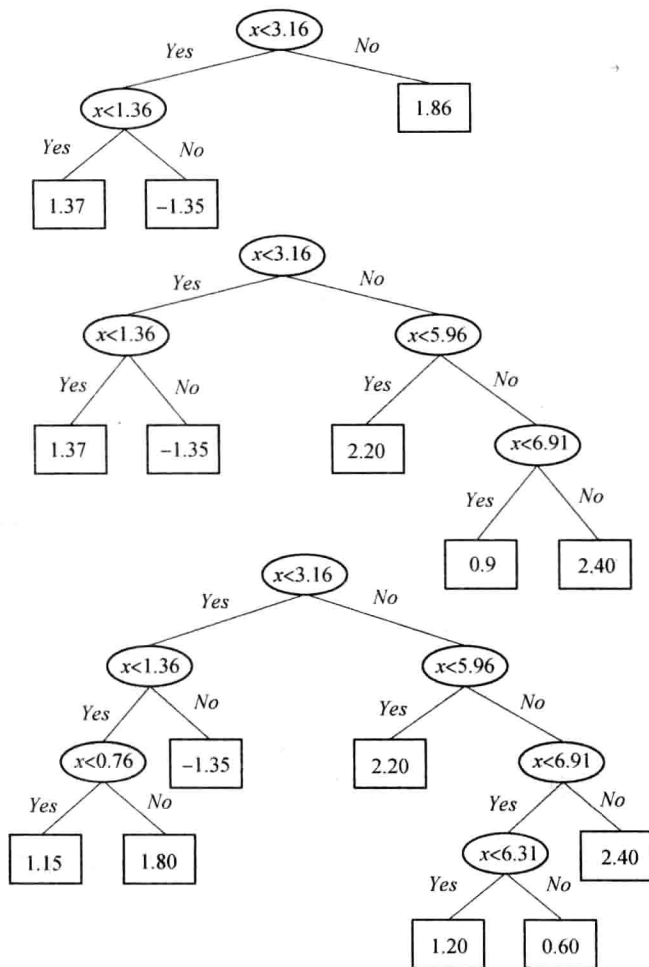


图 9-5 对于 θ 的不同值, 实现图 9-4 的光滑的回归树

195 先剪枝与后剪枝相比, 先剪枝较快, 但是后剪枝通常导致更准确的树。

9.4 由决策树提取规则

决策树能够提取特征。单变量树只使用必要的变量, 并且在树构建之后某些特征可能根本没有使用。我们还可以认为越靠近树根的特征从全局上讲越重要。例如, 图 9-6 中的决策树使用了变量 x_1 , x_2 和 x_4 , 但没有使用 x_3 。可以使用决策树提取特征: 构建一棵决策树, 并取该树使用的特征作为另一种学习方法的输入。

决策树的另一优点是可解释性(interpretability): 决策树节点中的条件简单、易于理解。从树根到树叶的每条路径对应于条件的合取, 这是因为到达树叶, 所有这些条件都必须满足。这些路径可以用 IF-THEN 规则集表示, 称作规则库(rule base)。一种这样的方法是

C4.5 规则(Quinlan 1993)。

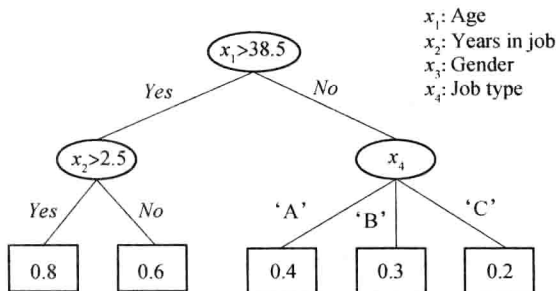


图 9-6 一棵(假想的)决策树。由根到树叶的每条路径都可以用一个合取规则表示, 由该路径上决策节点定义的条件组成

例如, 图 9-6 的决策树可以用如下规则集表示:

R1: IF(age > 38.5) AND(years-in-job > 2.5) THEN $y = 0.8$

R2: IF(age > 38.5) AND(years-in-job ≤ 2.5) THEN $y = 0.6$

R3: IF(age ≤ 38.5) AND(job-type = 'A') THEN $y = 0.4$

R4: IF(age ≤ 38.5) AND(job-type = 'B') THEN $y = 0.3$

R5: IF(age ≤ 38.5) AND(job-type = 'C') THEN $y = 0.2$

这样的规则库可以提取知识; 它容易理解, 并且使得领域专家可以验证从数据学习得到的模型。对于每个规则, 我们可以计算被该规则覆盖的训练数据所占的百分比, 即规则的支持度(rule support)。这些规则反映数据集的主要特性: 它们显示了重要特征和划分位置。例如, 在这个(假想的)例子中, 我们看到就我们的目的(y)而言, 38 岁或更年轻的人不同于 39 岁或更年长的人。并且, 在后一组, 工作类型区分他们; 而在前一组, 做一项工作的年限是最好的区分特征。

对于分类树, 可能有多个树叶被标记为相同的类。在这种情况下, 对应不同路径的多个合取表达式可以合并成一个析取(OR)。类区域对应于多个小区域的并, 而每个小区域对应一个树叶定义的区域。例如, 图 9-1 的 C_1 类可以表示为:

$$\text{IF}(x \leq w_{10}) \text{OR}((x_1 > w_{10}) \text{AND}(x_2 \leq w_{20})) \text{THEN } C_1$$

为了简化, 可以修剪规则(pruning rule)。剪掉一棵子树对应同时从一些规则剪去一些项。可以从一个规则剪去一个项而不涉及其他规则。例如, 在前面的规则集中, 对于 R3, 如果所有 job-type = 'A' 的人无论他的年龄多大, 都具有大致为 0.4 的输出, 则可以对 R3 剪枝, 得到

$$\text{R3': IF}(\text{job-type} = 'A') \text{ THEN } y = 0.4$$

注意, 规则剪枝后可能不能再写回到树中。

9.5 由数据学习规则

正如我们刚刚看到的, 产生 IF-THEN 规则的一种方法是训练一棵决策树, 并把它转换成规则; 另一种方法是直接学习规则。规则归纳(rule induction)类似于树归纳, 唯一的区别在于规则归纳进行深度优先搜索, 并且一次产生一条路径(规则); 而树归纳进行宽度优先搜索, 并且同时产生所有路径。

一次学习一个规则。每个规则是离散或数值属性上条件的合取(与决策树一样), 并且这些条件一次添加一个, 以优化某个标准, 如最小化熵。我们说规则覆盖(cover)一个实例, 如

果该实例满足规则的所有条件。一旦规则形成并被剪枝,就将它添加到规则库中,从训练集中删除被它覆盖的所有训练实例,并且继续该过程,直到得到足够的规则。这称作顺序覆盖(sequential covering)。外循环一次将一个规则添加到规则库,而内循环一次将一个条件添加到当前规则中。这些步骤都是贪心的,并且不能保证最优。为了更好地泛化,两个循环都有剪枝步。

规则归纳算法的一个例子是 *Ripper* (Cohen 1995), 它基于较早的算法 *Irep* (Fürnkranz 和 Widmer 1994)。我们从两类问题开始,并使用术语正例和负例,然后再推广到 $K > 2$ 类。添加规则旨在解释正例,使得如果一个实例不被任何规则覆盖,则它将被归到负类。这样,当规则匹配时,它或者是正确的(真正),或者导致一个假正。*Ripper* 的外循环的伪代码在图 9-7 中。

```

Ripper( Pos, Neg, k)
  RuleSet ← LearnRuleSet( Pos, Neg)
  For k times
    RuleSet ← OptimizeRuleSet( RuleSet, Pos, Neg)
  LearnRuleSet( Pos, Neg)
    RuleSet ← ∅
    DL ← DescLen( RuleSet, Pos, Neg)
    Repeat
      Rule ← LearnRule( Pos, Neg)
      将 Rule 添加到 RuleSet
      DL' ← DescLen( RuleSet, Pos, Neg)
      If DL' > DL + 64
        PruneRuleSet( RuleSet, Pos, Neg)
        Return RuleSet
      If DL' < DL DL ← DL'
      从 Pos 和 Neg 删除被 Rule 覆盖的实例
    Until Pos = ∅
  Return RuleSet
PruneRuleSet( RuleSet, Pos, Neg)
  For 每个 Rule ∈ RuleSet, 按相反次序
    DL ← DescLen( RuleSet, Pos, Neg)
    DL' ← DescLen( RuleSet - Rule, Pos, Neg)
    If DL' < DL 从 RuleSet 中删除 Rule
  Return RuleSet
OptimizeRuleSet( RuleSet, Pos, Neg)
  For 每个 Rule ∈ RuleSet
    DL0 ← DescLen( RuleSet, Pos, Neg)
    DL1 ← DescLen( RuleSet - Rule +
      ReplaceRule( RuleSet, Pos, Neg), Pos, Neg)
    DL2 ← DescLen( RuleSet - Rule +
      ReviseRule( RuleSet, Rule, Pos, Neg), Pos, Neg)
    If DL1 = min( DL0, DL1, DL2)
      从 RuleSet 中删除 Rule 并且
      添加 ReplaceRule( RuleSet, Pos, Neg)
    Else If DL2 = min( DL0, DL1, DL2)
      从 RuleSet 中删除 Rule 并且
      添加 ReviseRule( RuleSet, Rule, Pos, Neg)
  Return RuleSet

```

图 9-7 学习规则的 *Ripper* 算法。只给出了外循环,内循环与在决策树中添加一个节点类似

在 Ripper 中, 条件被添加到规则中, 以最大化 Quinlan 的 Foil 算法(1990)使用的信息增益度量。假设我们有规则 R , 并且 R' 是添加一个条件后的候选规则。增益的改变定义为

$$Gain(R', R) = s \cdot \left(\log_2 \frac{N'_+}{N'} - \log_2 \frac{N_+}{N} \right) \quad (9.17)$$

其中 N 是被 R 覆盖的实例数, 而 N_+ 是其中的真正例数。类似地, N' 是被 R' 覆盖的实例数, N'_+ 是其中的真正例数。 s 是 R 中的真正例并且增加条件之后在 R' 也是真正实例的实例数。根据信息理论, 增益的变化度量对一个正例编码所需位的减少。

向规则增加条件直到它不再覆盖负例。一旦规则形成, 就以相反的次序对它剪枝, 以便找到最大化规则价值度量(rule value metric)的规则

$$rvm(R) = \frac{p - n}{p + n} \quad (9.18)$$

其中 p 和 n 分别是剪枝集中的真正例和假正例数。剪枝集是数据的三分之一, 已经使用三分之二的数据作为增长集。

一旦规则形成并被剪枝, 就从训练集中删除被规则覆盖的所有正的和负的训练实例。如果还有正的实例, 则继续进行规则归纳。在存在噪声的情况下, 即当规则不能解释足够多的实例时, 我们可能提前中止归纳。为了度量规则的价值, 使用最小描述长度(参见 4.8 节)(Quinlan 1995)。典型地, 如果规则的描述长度不短于它所解释的实例的描述长度, 则我们停止。规则库的描述长度是规则库中所有规则的描述长度之和, 加上不被规则库覆盖的实例的描述长度。当规则的描述长度比迄今得到的最佳描述长度多 64 位时, Ripper 停止添加规则。一旦学到了规则库, 我们就以逆序忽略规则, 看是否能够删除它们而不增加描述长度。

规则库中的规则在学习之后也要优化。对一个规则, Ripper 考虑两种可供选择的方案: 一种是置换规则, 从空规则开始, 增长然后剪枝。第二种是修订规则, 从规则开始, 增长然后剪枝。这两个规则与原规则比较, 并将三个中的最短者添加到规则库中。规则库的这种优化进行 k 次, 通常进行两次。

当存在 $K > 2$ 个类时, 将这些类按照它们的先验概率排序, 使得 C_1 的先验概率最低, C_K 的先验概率最高。然后定义一系列两类问题。开始, 属于 C_1 的实例为正例, 其他类的实例都是负例。学习 C_1 的规则后, 删除它的所有实例, 学习将 C_2 与 C_3, \dots, C_K 分离开来。重复该过程, 直到只剩下 C_K 。空的缺省规则标记为 C_K , 使得如果一个实例不被任何规则覆盖, 则将它指派到 C_K 。

对于大小为 N 的训练集, Ripper 的复杂度为 $\mathcal{O}(N \log^2 N)$, 并且可以用于很大的训练集(Dietterich 1997)。学习的规则是命题规则(propositional rule)。更准确地说, 是条件中包含变量的一阶规则(first-order rule), 称作谓词(predicate)。谓词是一个函数, 依赖于其变元的值, 它返回真或假。因此, 谓词可以定义属性值之间的关系, 而命题不能(Mitchell, 1997):

IF Father(y, x) AND Female(y) THEN Daughter(x, y)

在逻辑程序设计语言(如 Prolog)中, 这种规则可以看作程序, 而从数据中学习它们称作归纳逻辑程序设计(inductive logic programming)。一种这样的算法是 Foil(Quinlan 1990)。

将一个值指派到一个变量称作绑定(binding)。如果训练集中存在到变量的绑定集, 则

称规则匹配。学习一阶规则类似于学习命题规则，外循环添加规则，而内循环相规则添加条件，在内循环结束时进行规则剪枝。不同的是，在内循环中，每一步我们考虑增加一个谓词（而不是命题），并检查规则的性能提高（Mitchell 1997）。为了评估规则的性能，我们考虑变量的所有可能绑定，对训练集中正的和负的绑定计数，并使用，例如式(9.17)。在学习一阶规则时，我们使用谓词而不是命题，因此这些谓词应当事先定义，并且训练集是已知为真的谓词集。

9.6 多变量树

在构造单变量树时，划分时只使用一个输入维。在构造多变量树(multivariate tree)时，在每个决策节点都可以使用所有的输入维，因此更加一般。当所有的输入都是数值属性时，二元线性多变量节点定义为

$$f_m(\mathbf{x}) : \mathbf{w}_m^T \mathbf{x} + w_{m0} > 0 \quad (9.19)$$

因为线性多变量节点取加权和，因此离散属性应当用0/1哑数值变量表示。式(9.19)定义了一个具有任意方向的超平面(参见图9-8)。从根到树叶的路径上的相继节点进一步划分实例，而叶节点定义输入空间上的多面体。具有数值特征的一元节点是一种特例，所有的 w_{mj} 除一个之外均为0。这样，式(9.1)的单变量数值节点也定义了一个线性判别式，但是与轴 x_j 正交于 w_{m0} ，与其他轴 x_i 平行。因此，我们看到在单变量节点有 d 个可能的方向(\mathbf{w}_m)和 $N_m - 1$ 个可能的阈值($-w_{m0}$)，使得穷举搜索是可能的。在多变量节点，有 $2^d C_{N_m}^d$ 个可能的超平面(Murthy、Kasif 和 Salzberg 1994)，并且不再可能进行穷举搜索。

202

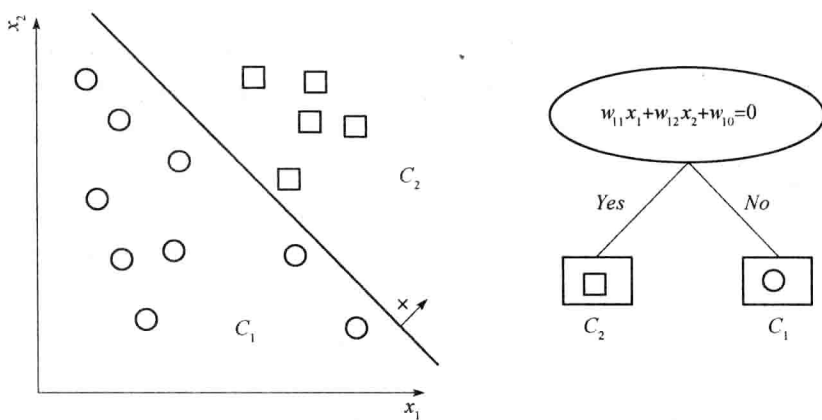


图9-8 线性多变量决策树。线性多变量节点可以安放任意超平面，因而更一般，而单变量节点限于平行于轴的划分

当我们从单变量节点过渡到线性多变量节点时，节点变得更灵活。使用非线性多变量节点，还可以更加灵活。例如，使用二次方程，我们有

$$f_m(\mathbf{x}) : \mathbf{x}^T \mathbf{W}_m \mathbf{x} + \mathbf{w}_m^T \mathbf{x} + w_{m0} > 0 \quad (9.20)$$

Guo 和 Gelfand(1992)提出使用多层感知器(见第11章)。多层感知器是非线性基函数的线性和，是另一种产生非线性决策节点的方法。还有一种可能性是使用球形节点(sphere node)

(Devroye、Györfi 和 Lugosi 1996)

$$f_m(\mathbf{x}) : \|\mathbf{x} - \mathbf{c}_m\| \leq \alpha_m \quad (9.21)$$

其中 \mathbf{c}_m 是球心, α_m 是半径。

已经提出了一些学习用于分类的多变量决策树的算法: 最早的算法是 CART 算法的多变量版本 (Breiman 等 1984), 它逐一对权重 w_{mj} 进行微调来降低不纯度。CART 还包含一个预处理步骤, 通过子集选择降低维度 (见第 6 章) 并降低节点的复杂度。一种对 CART 扩展的算法是 OC1 算法 (Murthy、Kasif 和 Salzberg 1994)。一种可能的方法 (Loh 和 Vanichsetakul 1988) 是假设所有的类是高斯的, 具有共同的协方差矩阵, 因此具有分离每个类的线性判别式 (第 5 章)。在这种情况下, 具有 K 个类, 每个节点具有 K 个分支, 而每个分支携带一个将每个类与其他类分开的线性判别式。Brodley 和 Utgof (1995) 提出了一种方法, 这种方法训练线性判别式以最小化分类误差 (第 10 章)。Guo 和 Gelfand (1992) 提出了一种方法, 将 $K > 2$ 个类划分为两个超群, 然后学习二元多变量树。Loh 和 Shih (1997) 使用 2-均值聚类 (第 7 章) 将数据聚成两组。一旦类聚成两组, Yıldız 和 Alpaydm (2000) 就使用 LDA (第 6 章) 找出判别式。

203

任何分类器都近似一个从假设类中选取一个假设的实际 (未知的) 判别式。当我们使用单变量节点时, 我们的近似使用分段的、平行于轴的超平面。使用线性多变量节点, 我们可以使用任意的超平面, 并且使用较少的节点得到更好的近似。如果潜在的判别式是曲的, 非线性节点更好。分支因子确定节点定义的判别式的个数, 具有类似效果。具有两个分支的二元决策节点定义一个将输入空间一分为二的判别式。 n -路节点将输入空间划分成 n 部分。这样, 节点、分支因子和树大小之间存在相关性。使用简单节点和较低的分支因子可以得到一棵大树。但是, 这样的树 (如具有单变量二元节点) 可解释性更好。线性多变量节点更难解释。更复杂的节点也需要更多的数据, 并且随着我们沿树向下, 数据越来越少, 更容易过分拟合。如果节点复杂并且树比较小, 那么我们就失去了通过树想要得到的主要东西——将问题划分成一系列简单问题。毕竟, 我们可以在根节点具有一个非常复杂的分类器, 它区分所有的类, 但这样的话, 这不是一棵树了!

9.7 注释

自从凯撒将一个复杂的问题 (如高卢人问题) 分解成一组较简单的问题以来, 分治一直被作为一种启发式方法频繁使用。在计算机科学中, 频繁地使用树将复杂度从线性降低到对数时间。Breiman 等 1984 使得决策树在统计学流行, Quinlan 1986、1993 使得决策树在在机器学习中很流行。多变量树归纳方法最近才开始流行, Yıldız 和 Alpaydm 2000 给出了综述和许多数据集上的比较。许多研究者 (如 Guo 和 Gelfand 1992) 将树的简单性与多层感知器的准确性结合在一起 (见第 11 章)。然而, 许多研究表明单变量树相当准确、具有很好的可解释性, 并且线性 (非线性) 多变量节点带来的附加的复杂度很难被认为是合理的。

204

杂变量决策树 (omnivariate decision tree) (Yıldız 和 Alpaydm 2001) 是一种混合树结构, 其中树可以具有单变量、线性多变量和非线性多变量节点。其基本思想是在树构造期间, 每个决策节点对应于到达该节点的训练数据子集定义的一个不同的子问题, 不同的模型可能更合适, 应当找出和使用合适的模型。到处都用相同类型的节点相当于假定输入空间的所有部分

都具有相同的归纳偏倚。在杂变量树中,在每个节点,训练并在确认集上使用统计检验(第14章)比较不同类型的候选节点,确定哪一个泛化性能最好。除非复杂的决策节点表现出显著较高的准确率,否则就选取较简单的候选节点。结果表明,在树构造的早期更靠近树根的地方使用较复杂的节点,并且随着我们沿树向下,简单的单变量节点就足够了。随着我们越来越靠近树叶,问题越来越简单,同时,数据越来越少。在这种情况下,复杂的节点过分拟合,并且被统计检验拒绝。随着我们沿树向下,节点的个数指数增加。因此,大部分节点是单变量的,并且总体复杂度增加不太多。

决策树更多地用于分类而不是回归。它们非常流行:它们的学习和响应速度都很快,并且在许多领域都很准确(Murthy 1998)。由于它们的可解释性,甚至在有更准确的方法时,决策树仍是首选。当决策树写成 IF-THEN 规则集时,树可以理解,并且可以被具有应用领域知识的专家验证。

通常,在使用更复杂算法之前,建议先试验决策树,并将它的准确率作为性能基准。树分析还能帮助我们理解重要特征,单变量树还可以用于自动特征提取。单变量树的另一个重要优点是它可以使用数值和离散特征,而不需要将一种类型转换成另一种。

决策树是非参数方法,类似于第8章讨论的方法,但是存在一些不同之处:

- 每个树叶对应于一个“箱”,只不过箱不必具有相同的大小(如 Parzen 窗口)或相同个数的实例(如 k -最近邻)。
- 箱的划分不仅仅根据输入空间中的相似度,而是需要通过熵或均方误差,使用输出信息。
- 决策树的另一个优点是仅通过少量比较就能找到树叶(箱)。
- 决策树一旦构造就不需要存放所有的训练集,而只需要存放树的结构、决策节点的参数和树叶节点的输出值。与需要存储所有训练实例的基于核的或基于近邻的方法相比,这意味决策树的空间复杂度也非常小。

使用决策树,一个类不必具有所有实例都匹配的单个描述。一个类可以具有多个可能的描述,它们甚至可能在输入空间中不相交。

树不同于上一章讨论的统计模型。树直接地对分离类实例的判别式编码,而不必携带许多关于这些实例在该区域中如何分布的信息。决策树是基于判别式的(discriminant-based),而统计学方法是基于似然的(likelihood-based),它们在使用贝叶斯规则和计算判别式之前显式地估计 $p(\mathbf{x} | C_i)$ 。基于判别式的方法绕过类密度估计,直接估计判别式。在其后几章,我们将进一步讨论基于判别式的方法。

9.8 习题

1. 将基尼指数式(9.5)和误分类误差式(9.6)推广到 $K > 2$ 个类。考虑损失函数,将误分类误差推广到风险。
2. 对于数值属性,我们可以不用二元划分,而通过两个阈值和三个分支使用三元划分

$$x_j < w_{ma}, w_{ma} \leq x_j < w_{mb}, x_j \geq w_{mb}$$

修改决策树归纳方法,学习两个阈值 w_{ma} 和 w_{mb} 。与二元节点相比,这种节点的优缺点是什么?

3. 提出一种带回溯的树归纳算法。
4. 在产生单变量树时, 具有 n 个可能值的离散属性可以用 n 个 0/1 哑变量表示, 并将它们看作是 n 个分别的数值属性。这种方法的优缺点是什么?
5. 为球形树式(9.21)推导一个学习算法。将它推广到椭球形树。
6. 在回归树中, 我们提到在树叶节点不是计算均值, 而是可以做线性回归拟合, 并使树叶上的响应依赖于输入。对分类树提出一种类似的方法。
7. 为回归提出一种规则归纳算法。
8. 在回归树中, 如何消除树叶边界的不连续性?
9. 假设对于分类问题我们已经有一棵训练后的决策树, 那么除训练集之外, 如何在构建 k 最近邻分类时使用它?
10. 在多变量树中, 很可能在一个节点上不需要所有的输入变量。如何降低节点的维度?

9.9 参考文献

- Breiman, L., J. H. Friedman, R. A. Olshen, and C. J. Stone. 1984. *Classification and Regression Trees*. Belmont, CA: Wadsworth International Group.
- Brodley, C. E., and P. E. Utgoff. 1995. "Multivariate Decision Trees." *Machine Learning* 19: 45-77.
- Cohen, W. 1995. "Fast Effective Rule Induction." In *Twelfth International Conference on Machine Learning*, ed. A. Frieditis and S. J. Russell, 115-123. San Mateo, CA: Morgan Kaufmann.
- Devroye, L., L. Györfi, and G. Lugosi. 1996. *A Probabilistic Theory of Pattern Recognition*. New York: Springer.
- Dietterich, T. G. 1997. "Machine Learning Research: Four Current Directions." *AI Magazine* 18: 97-136.
- Fürnkranz, J., and G. Widmer. 1994. "Incremental Reduced Error Pruning." In *Eleventh International Conference on Machine Learning*, ed. W. Cohen and H. Hirsh, 70-77. San Mateo, CA: Morgan Kaufmann.
- Guo, H., and S. B. Gelfand. 1992. "Classification Trees with Neural Network Feature Extraction." *IEEE Transactions on Neural Networks* 3: 923-933.
- Loh, W.-Y., and Y. S. Shih. 1997. "Split Selection Methods for Classification Trees." *Statistica Sinica* 7: 815-840.
- Loh, W.-Y., and N. Vanichsetakul. 1988. "Tree-Structured Classification via Generalized Discriminant Analysis." *Journal of the American Statistical Association* 83: 715-725.
- Mitchell, T. 1997. *Machine Learning*. New York: McGraw-Hill.
- Murthy, S. K. 1998. "Automatic Construction of Decision Trees from Data: A Multi-Disciplinary Survey." *Data Mining and Knowledge Discovery* 4: 345-389.
- Murthy, S. K., S. Kasif, and S. Salzberg. 1994. "A System for Induction of Oblique Decision Trees." *Journal of Artificial Intelligence Research* 2: 1-32.
- Quinlan, J. R. 1986. "Induction of Decision Trees." *Machine Learning* 1: 81-106.
- Quinlan, J. R. 1990. "Learning Logical Definitions from Relations." *Machine Learning* 5: 239-266.

- Quinlan, J. R. 1993. *C4.5: Programs for Machine Learning*. San Mateo, CA: Morgan Kaufmann.
- Quinlan, J. R. 1995. "MDL and Categorical Theories (continued)." In *Twelfth International Conference on Machine Learning*, ed. A. Prieditis and S. J. Russell, 467-470. San Mateo, CA: Morgan Kaufmann.
- Rokach, L., and O. Maimon. 2005. "Top-Down Induction of Decision Trees Classifiers—A Survey." *IEEE Transactions on Systems, Man, and Cybernetics—Part C* 35: 476-487.
- Yıldız, O. T., and E. Alpaydın. 2000. "Linear Discriminant Trees." In *Seventeenth International Conference on Machine Learning*, ed. P. Langley, 1175-1182. San Francisco: Morgan Kaufmann.
- Yıldız, O. T., and E. Alpaydın. 2001. "Omnivariate Decision Trees." *IEEE Transactions on Neural Networks* 12: 1539-1546.

线性判别式

在线性判别式中, 我们假定类的实例是线性可分的。这是一种基于判别式的方法, 它直接估计判别式的参数, 而不必先估计各种概率。

10.1 引言

在前面的章节中, 对于分类, 我们定义了一组判别式函数 $g_j(\mathbf{x})$, $j = 1, \dots, K$, 并且如果 $g_j(\mathbf{x}) = \max_{j=1}^K g_j(\mathbf{x})$, 我们就选择 C_i 。

前面, 在我们讨论分类方法时, 我们首先估计先验概率 $\hat{p}(C_i)$ 和类似然 $\hat{p}(\mathbf{x} | C_i)$, 再使用贝叶斯规则计算后验密度。然后, 我们使用后验密度定义判别式函数, 例如

$$g_i(\mathbf{x}) = \log \hat{p}(C_i | \mathbf{x})$$

这称作基于似然的分类 (likelihood-based classification), 并且我们在前面已经讨论了估计类似然 $P(\mathbf{x} | C_i)$ 的参数 (第 5 章)、半参数 (第 7 章) 和非参数 (第 8 章) 方法。

现在, 我们讨论基于判别式的分类 (discriminant-based classification), 这里我们绕过似然或后验概率的估计, 直接为判别式假定模型。基于判别式的方法对类之间的判别式形式进行假设, 而不对密度 (例如, 是否是高斯分布)、输入是否等相关知识做任何假设。

209

我们为判别式定义一个模型

$$g_i(\mathbf{x} | \Phi_i)$$

显式地用参数 Φ_i 的集合参数化。这与基于似然的模式不同。基于似然的方法在定义似然密度时具有隐式参数。这是不同的归纳偏倚: 我们对判别式的形式进行假设, 而不是对密度的形式进行假设。

学习是优化模型参数 Φ_i , 最大化给定类标号的训练集上的分类准确率。这不同于基于似然的方法。基于似然的方法分别为每个类搜索最大化样本似然的参数。

在基于判别式的方法中, 我们并不关注正确地估计类区域中的密度; 我们所关注的是正确估计类区域之间的边界 (boundary)。基于判别式方法的创导者 (如 Cherkassky 和 Mulier 1998) 指出, 估计类密度比估计类判别式更困难, 并且为解决较容易的问题而解决困难的问题并无意义。当然, 仅当判别式可以用简单函数近似时才确实如此。

本章, 我们关注最简单的情况, 其中判别式是 \mathbf{x} 的线性函数:

$$g_i(\mathbf{x} | \mathbf{w}_i, w_{i0}) = \mathbf{w}_i^T \mathbf{x} + w_{i0} = \sum_{j=1}^d w_{ij} x_j + w_{i0} \quad (10.1)$$

线性判别式 (linear discriminant) 经常使用, 主要是由于它的简单性, 即它的空间和时间复杂度都是 $\mathcal{O}(d)$ 。线性模型容易理解: 最终的输出是若干因素的加权和。权重的大小显示

了这些因素的重要性, 而它们的符号显示其作用的正负。大部分函数是可加的, 因为输出是若干属性作用的加权和, 其中权重可能是正的(加强)或负的(抑制)。例如, 当一位顾客申请信用卡时, 金融机构计算申请者的信用得分。得分一般是多个属性作用之和。例如, 年薪的作用为正(较高的年薪增加得分)。

在许多应用中, 线性判别式相当准确。例如, 我们知道当类是高斯的, 具有相同的协方差矩阵时, 最佳的判别式是线性的。然而, 即使该假设不成立, 也可以使用线性判别式, 并且不必对类密度做任何假设就能计算模型参数。在试用更复杂的模型, 确保附加的复杂性是合理的之前, 我们将一直使用线性判别式。

正如我们一直做的那样, 我们把寻找线性判别式函数问题归结为搜索最小化某个误差函数的参数值问题。我们尤其关注优化准则函数的梯度(gradient)方法。

10.2 推广线性模型

当线性判别式不够灵活时, 我们可以提高复杂度, 使用二次判别式(quadratic discriminant)函数

$$g_i(\mathbf{x} | \mathbf{W}_i, \mathbf{w}_i, w_{i0}) = \mathbf{x}^T \mathbf{W}_i \mathbf{x} + \mathbf{w}_i^T \mathbf{x} + w_{i0} \quad (10.2)$$

但是, 这种方法的复杂度是 $\mathcal{O}(d^2)$, 并且我们还会遇到偏倚和方差的两难选择: 尽管二次模型更一般, 但是它需要更大的训练集, 并且在小样本上可能过分拟合。

一种等价的方法是通过增加高阶项(higher-order term), 又称乘积项(product term), 对输入进行预处理。例如, 对于两个输入 x_1 和 x_2 , 我们可以定义新变量

$$z_1 = x_1, z_2 = x_2, z_3 = x_1^2, z_4 = x_2^2, z_5 = x_1 x_2$$

并取 $\mathbf{z} = [z_1, z_2, z_3, z_4, z_5]^T$ 为输入。定义在五维 \mathbf{z} 空间上的线性函数对应二维 \mathbf{x} 空间上的非线性函数。替代在原空间定义非线性函数(判别式或回归), 我们需要做的是定义到新空间的、合适的非线性变换, 其中新空间上的函数可以是线性的。

判别式可以表示成

$$g_i(\mathbf{x}) = \sum_{j=1}^k w_j \phi_{ij}(\mathbf{x}) \quad (10.3)$$

其中 $\phi_{ij}(\mathbf{x})$ 是基函数(basis function)。例子如下

- $\sin(x_1)$
- $\exp(-(x_1 - m)^2/c)$
- $\exp(-\|\mathbf{x} - \mathbf{m}\|^2/c)$
- $\log(x_2)$
- $1(x_1 > c)$
- $1(ax_1 + bx_2 > c)$

其中 m, a, b, c 是标量, \mathbf{m} 是 d 维向量, 而当 b 为真时 $1(b)$ 返回1, 否则返回0。将非线性函数表示成非线性基函数的线性和的想法并非新想法, 并且最初称作势函数(potential function)(Aizerman, Braverman 和 Rozonoer 1964)。多层感知器(见第11章)和径向基函数(见第12章)具有进一步的优点, 可以在学习时调整基函数的参数。13章我们讨论支持向量机, 它使用由这种基函数构造的核函数。

10.3 线性判别式的几何意义

10.3.1 两类问题

让我们从最简单的两类问题开始。在这种情况下，一个判别式函数就足够了：

$$\begin{aligned} g(\mathbf{x}) &= g_1(\mathbf{x}) - g_2(\mathbf{x}) \\ &= (\mathbf{w}_1^T \mathbf{x} + w_{10}) - (\mathbf{w}_2^T \mathbf{x} + w_{20}) \\ &= (\mathbf{w}_1 - \mathbf{w}_2)^T \mathbf{x} + (w_{10} - w_{20}) \\ &= \mathbf{w}^T \mathbf{x} + w_0 \end{aligned}$$

并且如果 $g(\mathbf{x}) > 0$ ，我们选择 C_1 ，否则选择 C_2 。

这定义了一个超平面，其中 \mathbf{w} 是权重向量 (weight vector)， w_0 是阈值 (threshold)。后者称作阈值是因为规则可以改写为：如果 $\mathbf{w}^T \mathbf{x} > -w_0$ ，选择 C_1 ，否则选择 C_2 。超平面将输入空间划分成两个半空间： C_1 的决策区域 R_1 和 C_2 的决策区域 R_2 。 R_1 中的任何 \mathbf{x} 都在超平面的正 (positive) 侧，而 R_2 中的任何 \mathbf{x} 都在超平面的负 (negative) 侧。当 \mathbf{x} 为 0 时， $g(\mathbf{x}) = w_0$ 并且如果我们有 $w_0 > 0$ ，则原点在超平面的正侧，如果 $w_0 < 0$ ，则原点在超平面的负侧，而如果 $w_0 = 0$ ，则超平面经过原点 (见图 10-1)。

取决策面上的两个点 \mathbf{x}_1 和 \mathbf{x}_2 (即 $g(\mathbf{x}_1) = g(\mathbf{x}_2) = 0$)，则

$$\mathbf{w}^T \mathbf{x}_1 + w_0 = \mathbf{w}^T \mathbf{x}_2 + w_0$$

$$\mathbf{w}^T (\mathbf{x}_1 - \mathbf{x}_2) = 0$$

并且我们看到 \mathbf{w} 是超平面上的任意向量的法线。让我们将 \mathbf{x} 改写为 (Duda、Hart 和 Stork 2001)

$$\mathbf{x} = \mathbf{x}_p + r \frac{\mathbf{w}}{\|\mathbf{w}\|}$$

其中 \mathbf{x}_p 是 \mathbf{x} 到超平面的法向投影，而 r 给出 \mathbf{x} 到超平面的距离，如果 \mathbf{x} 在负侧，则它为负；如果 \mathbf{x} 在正侧，则它为正 (参见图 10-2)。计算 $g(\mathbf{x})$ 并注意 $g(\mathbf{x}_p) = 0$ ，我们有

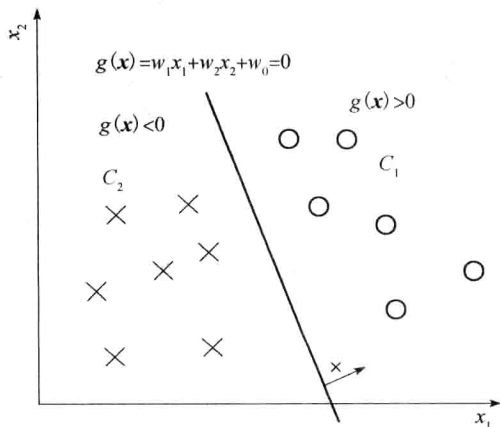


图 10-1 在二维情况下，线性判别式是一条将两个类的实例分开的直线

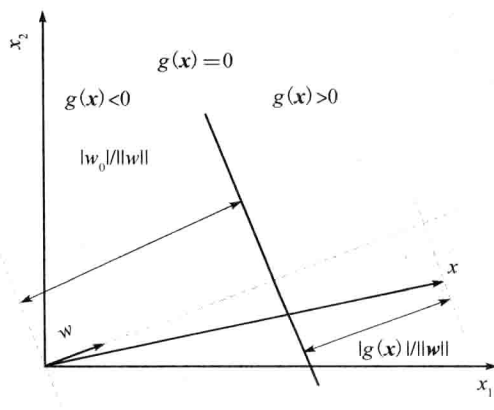


图 10-2 线性判别式的几何解释

$$r = \frac{g(\mathbf{x})}{\|\mathbf{w}\|} \quad (10.4)$$

于是, 我们看到超平面到原点的距离为

$$r_0 = \frac{w_0}{\|\mathbf{w}\|} \quad (10.5)$$

这样, w_0 决定超平面关于原点的位置, 而 \mathbf{w} 决定它的方向。

10.3.2 多类问题

当存在 $K > 2$ 个类时, 有 K 个判别式函数。当它们都是线性的时, 我们有

$$g_i(\mathbf{x} | \mathbf{w}_i, w_{i0}) = \mathbf{w}_i^T \mathbf{x} + w_{i0} \quad (10.6)$$

稍后, 我们将讨论学习, 但是现在我们假定参数 \mathbf{w}_i 和 w_{i0} 的计算使得对于训练集中的所有 \mathbf{x} ,

$$g_i(\mathbf{x} | \mathbf{w}_i, w_{i0}) = \begin{cases} > 0 & \text{如果 } \mathbf{x} \in C_i \\ \leq 0 & \text{否则} \end{cases} \quad (10.7)$$

使用这种判别函数相当于假设所有的类都是线性可分的 (linearly separable); 即对于每个类 C_i , 存在一个超平面 H_i , 使得所有的 $\mathbf{x} \in C_i$ 都在它的正侧, 所有的 $\mathbf{x} \in C_j, j \neq i$ 都在它的负侧 (参见图 10-3)。

在检验时, 给定 \mathbf{x} , 理想情况下应当只有一个 $g_j(\mathbf{x}) (j=1, \dots, K)$ 大于 0, 而其他的都小于 0。但是, 并非总是如此: 这些超平面的正的半个空间可能重叠, 或者说可能存在所有的 $g_j(\mathbf{x})$ 都小于 0 的案例。这些案例可以看作拒绝 (reject) 案例, 但是通常的方法是将 \mathbf{x} 指派到具有最大判别式值的类:

$$\text{选择 } C_i \text{ 如果 } g_i(\mathbf{x}) = \max_{j=1}^K g_j(\mathbf{x}) \quad (10.8)$$

注意, $|g_i(\mathbf{x})| / \|\mathbf{w}_i\|$ 是从输入点 \mathbf{x} 到超平面的距离。假定 \mathbf{w}_i 具有类似的长度, 这将该点指派到这个类, (在所有的 $g_j(\mathbf{x}) > 0$ 中) 该点到其超平面最远。这称作线性分类器 (linear classifier), 并且它将特征空间几何地划分成 K 个凸决策区域 R_i (参见图 10-3)。

10.4 逐对分离

如果类不是线性可分的, 一种方法是将它分成一组线性问题。一种可能的方法是逐对分离 (pairwise separation) (Duda、Hart 和 Stork 2001)。它使用 $K(K-1)/2$ 个线性判别式 $g_{ij}(\mathbf{x})$, 每对不同的类一个 (参见图 10-4):

$$g_{ij}(\mathbf{x} | \mathbf{w}_{ij}, w_{ij0}) = \mathbf{w}_{ij}^T \mathbf{x} + w_{ij0}$$

参数 $w_{ij} (j \neq i)$ 在训练时计算, 使得

$$g_{ij}(\mathbf{x}) = \begin{cases} > 0 & \text{如果 } \mathbf{x} \in C_i \\ \leq 0 & \text{如果 } \mathbf{x} \in C_j \\ \text{无定义} & \text{否则} \end{cases} \quad i, j = 1, \dots, K \text{ 并且 } i \neq j \quad (10.9)$$

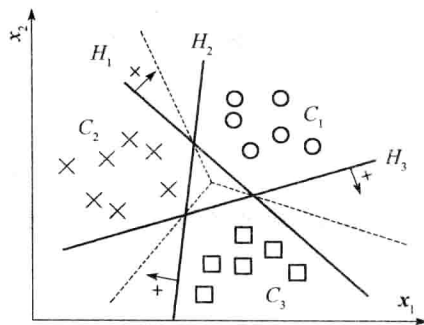


图 10-3 在线性分类, 每个超平面 H_i 将 C_i 类的实例与其他类的实例分开。为了做到这一点, 类应当是线性可分的。虚线是线性分类器的归约后的边界

也就是说, 如果 $\mathbf{x}^t \in C_k$, $k \neq i$, $k \neq j$, 则在 $g_{ij}(\mathbf{x})$ 训练时不使用 \mathbf{x}^t 。

在检验时, 如果对任意的 $j \neq i$, 都有 $g_{ij}(\mathbf{x}) > 0$, 则我们选择 C_i 。

在许多情况下, 可能对于任何 i , 这一条件都不满足, 而如果我们不想丢弃这种情况的话, 我们可以用如下和式放宽该合取, 选取最大化下式的类:

$$g_i(\mathbf{x}) = \sum_{j \neq i} g_{ij}(\mathbf{x}) \quad (10.10)$$

即使这些类不是线性可分的, 如果这些类是逐对线性可分的(这种情况可能性更大), 则可以使用逐对分离, 导致类的非线性分离(参见图 10-4)。这是将复杂问题(例如, 非线性问题)分解成一系列较简单问题(例如, 线性问题)的又一种方法。我们已经看到使用这一思想的决策树(第 9 章), 并且在第 17 章, 我们还将看到组合多个模型的更多例子, 例如, 交错输出码和混合专家模型, 其中线性模型数小于 $\mathcal{O}(K^2)$ 。

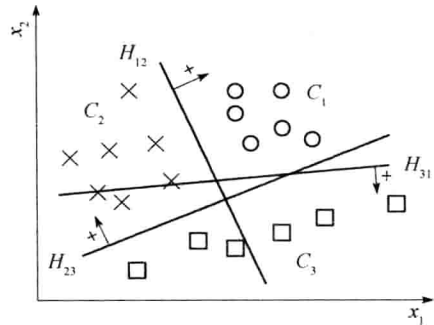


图 10-4 在逐对线性分离中, 每一对类有一个分离超平面。一个输入被指派到 C_1 , 它应当在 H_{12} 和 H_{13} 的正侧(H_{13} 的正侧是 H_{31} 的负侧); 我们不考虑 H_{23} 的值。在这种情况下, C_1 不是关于其他类线性可分的, 但却是逐对线性可分的

10.5 参数判别式的进一步讨论

在第 5 章, 我们看到如果类密度 $p(\mathbf{x} | C_i)$ 是高斯的, 并且具有共同的协方差矩阵, 则判别式函数是线性的

$$g_i(\mathbf{x}) = \mathbf{w}_i^T \mathbf{x} + w_{i0} \quad (10.11)$$

其中参数可以用下式解析地计算

$$\begin{aligned} \mathbf{w}_i &= \Sigma^{-1} \boldsymbol{\mu}_i \\ w_{i0} &= -\frac{1}{2} \boldsymbol{\mu}_i^T \Sigma^{-1} \boldsymbol{\mu}_i + \log P(C_i) \end{aligned} \quad (10.12)$$

给定数据集, 我们首先计算 $\boldsymbol{\mu}_i$ 和 Σ 的估计, 然后把估计 \mathbf{m}_i 和 \mathbf{S} 插入式(10.12), 并计算线性判别式的参数。

让我们再次考虑两类的特殊情况。我们定义 $y \equiv P(C_1 | \mathbf{x})$, $p(C_2 | \mathbf{x}) = 1 - y$ 。则在分类时, 我们

$$\text{选择 } C_1, \text{ 如果 } \begin{cases} y > 0.5 \\ \frac{y}{1-y} > 1, \\ \log \frac{y}{1-y} > 0 \end{cases} \quad \text{否则选择 } C_2。$$

$\log y / (1 - y)$ 称作分对数(logit)变换或 y 的对数几率(log odd)。在两个共享相同的协方

差矩阵的正态类的情况下，对数几率是线性的：

$$\begin{aligned}
 \text{logit}(P(C_1 | \mathbf{x})) &= \log \frac{P(C_1 | \mathbf{x})}{1 - P(C_1 | \mathbf{x})} = \log \frac{P(C_1 | \mathbf{x})}{P(C_2 | \mathbf{x})} \\
 &= \log \frac{p(\mathbf{x} | C_1)}{p(\mathbf{x} | C_2)} + \log \frac{P(C_1)}{P(C_2)} \\
 &= \log \frac{(2\pi)^{-d/2} |\Sigma|^{-1/2} \exp[-(1/2)(\mathbf{x} - \boldsymbol{\mu}_1)^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}_1)]}{(2\pi)^{-d/2} |\Sigma|^{-1/2} \exp[-(1/2)(\mathbf{x} - \boldsymbol{\mu}_2)^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}_2)]} + \log \frac{P(C_1)}{P(C_2)} \\
 &= \mathbf{w}^T \mathbf{x} + w_0
 \end{aligned} \tag{10.13}$$

其中

$$\begin{aligned}
 \mathbf{w} &= \Sigma^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) \\
 w_0 &= -\frac{1}{2} (\boldsymbol{\mu}_1 + \boldsymbol{\mu}_2)^T \Sigma^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) + \log \frac{P(C_1)}{P(C_2)}
 \end{aligned} \tag{10.14}$$

分对数的逆

$$\log \frac{P(C_1 | \mathbf{x})}{1 - p(C_1 | \mathbf{x})} = \mathbf{w}^T \mathbf{x} + w_0$$

是逻辑斯谛 (logistic) 函数，又称 S 形 (sigmoid) 函数 (见图 10-5)：

$$P(C_1 | \mathbf{x}) = \text{sigmoid}(\mathbf{w}^T \mathbf{x} + w_0) = \frac{1}{1 + \exp[-(\mathbf{w}^T \mathbf{x} + w_0)]} \tag{10.15}$$

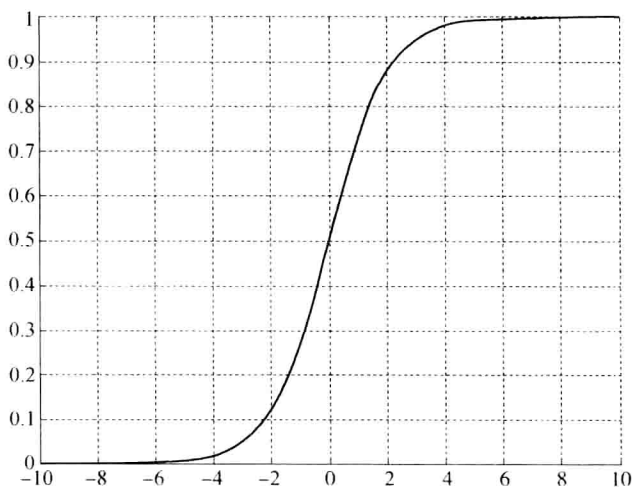


图 10-5 逻辑斯谛 (或 S 形) 函数

在训练阶段，我们估计 \mathbf{m}_1 , \mathbf{m}_2 , \mathbf{S} ，并将这些估计插入式 (10.14)，计算判别式的参数。在检验阶段，给定 \mathbf{x} ，我们可以

- 1) 计算 $g(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$ ，并且如果 $g(\mathbf{x}) > 0$ ，则选择 C_1 ；或者
- 2) 计算 $y = \text{sigmoid}(\mathbf{w}^T \mathbf{x} + w_0)$ ，并且如果 $y > 0.5$ ，则选择 C_1 ，因为 $\text{sigmoid}(0) = 0.5$ 。

在后一种情况下，S 形函数将判别式的值变换为后验概率。当有两个类并且只有一个判别式时，这是有效的。在 10.7 节，我们将讨论如何对 $K > 2$ 估计后验概率。

10.6 梯度下降

在基于似然的分类，参数是 $p(\mathbf{x} | C_i)$ 和 $P(C_i)$ 的有效统计量，而我们使用的估计参数的方法是最大化似然。在基于判别式的方法中，参数是判别式的参数，并且它们在最小化训练集上的分类误差是最优的。当 \mathbf{w} 表示参数集， $E(\mathbf{w} | \mathcal{X})$ 表示参数 \mathbf{w} 在给定训练集 \mathcal{X} 上的误差时，我们寻找

$$\mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{argmin}} E(\mathbf{w} | \mathcal{X})$$

在许多情况下，其中的一些稍后我们将看到，不存在解析解，而我们需要求助于迭代优化方法。最常用的方法是梯度下降 (gradient descent) 方法：当 $E(\mathbf{w})$ 是变量向量的可微函数时，我们有偏数组成的梯度向量 (gradient vector)

$$\nabla_{\mathbf{w}} E = \left[\frac{\partial E}{\partial w_1}, \frac{\partial E}{\partial w_2}, \dots, \frac{\partial E}{\partial w_d} \right]^T$$

和梯度下降 (gradient descent) 过程来最小化 E 。该方法从随机向量 \mathbf{w} 开始，并在每一步沿与该梯度相反的方向更新 \mathbf{w}

$$\Delta w_i = -\eta \frac{\partial E}{\partial w_i}, \forall i \quad (10.16)$$

$$w_i = w_i + \Delta w_i \quad (10.17)$$

其中 η 称作步长 (stepsize) 或学习因子 (learning factor)，决定向该方向移动多少。梯度上升用来最大化函数，并沿着梯度的方向前进。当我们得到极小 (或极大) 值时，导数等于 0，过程终止。这表明过程找到了一个最近的极小值，可能是局部极小值。除非函数只有一个极小值，否则不能保证找到全局极小。使用较好的 η 值也是至关重要的。如果太小，收敛可能太慢；太大可能导致摆动甚至发散。

在本书中，我们使用的梯度方法很简单，并且相当有效。然而，我们要记住，一旦确定了合适的方法和误差函数，就可以使用多种可能技术中的一种来优化模型参数，以便最小化误差函数。存在一些二阶方法和共轭梯度，收敛更快，但内存开销和计算量更大。像模拟退火和遗传算法这样的开销更大的方法可以更彻底地搜索参数空间，而不太依赖初始点的选择。

10.7 逻辑斯谛判别式

10.7.1 两类问题

在逻辑斯谛判别式 (logistic discrimination) 中，我们不是对类条件密度 $p(\mathbf{x} | C_i)$ ，而是对它们的比率建模。让我们还是从两类问题开始，并假定对数似然比是线性的：

$$\log \frac{p(\mathbf{x} | C_1)}{p(\mathbf{x} | C_2)} = \mathbf{w}^T \mathbf{x} + w_0 \quad (10.18)$$

当类条件密度为正态时 (参见式 (10.13))，这种假设确实成立。但是，逻辑斯谛判别式具有更广泛的应用。例如， \mathbf{x} 可能由离散属性组成，或者可能是连续和离散属性的混合。

使用贝叶斯规则，我们有

218
219

$$\begin{aligned}
 \text{logit}(P(C_1 | \mathbf{x})) &= \log \frac{P(C_1 | \mathbf{x})}{1 - P(C_1 | \mathbf{x})} \\
 &= \log \frac{p(\mathbf{x} | C_1)}{p(\mathbf{x} | C_2)} + \log \frac{P(C_1)}{P(C_2)} \\
 &= \mathbf{w}^T \mathbf{x} + w_0
 \end{aligned} \tag{10.19}$$

其中

$$w_0 = w_0^o + \log \frac{P(C_1)}{P(C_2)} \tag{10.20}$$

重新整理, 我们又得到 S 形函数

$$y = \hat{p}(C_1 | \mathbf{x}) = \frac{1}{1 + \exp[-(\mathbf{w}^T \mathbf{x} + w_0)]} \tag{10.21}$$

作为 $P(C_1 | \mathbf{x})$ 的估计。

让我们看看如何学习 \mathbf{w} 和 w_0 : 给定两个类的样本 $\mathcal{X} = \{\mathbf{x}^t, r^t\}$, 其中如果 $\mathbf{x} \in C_1$ 则 $r^t = 1$, 如果 $\mathbf{x} \in C_2$ 则 $r^t = 0$ 。我们假定给定 \mathbf{x}^t, r^t 是伯努利分布, 具有式(10.21)计算的概率 $y^t = P(C_1 | \mathbf{x}^t)$:

$$r^t | \mathbf{x}^t \sim \text{Bernoulli}(y^t)$$

这里, 我们看到了基于似然的方法与基于判别式的方法的区别: 对于前者, 我们对 $p(\mathbf{x} | C_i)$ 建模; 对于后者, 我们直接对 $r | \mathbf{x}$ 建模。样本的似然是

$$I(\mathbf{w}, w_0 | \mathcal{X}) = \prod_t (y^t)^{(r^t)} (1 - y^t)^{(1-r^t)} \tag{10.22}$$

我们知道, 当我们有一个需要最大化的似然函数时, 我们总是将它转换成需要最小化的误差函数 $E = -\log I$, 并且在我们的问题中, 我们有互熵 (cross-entropy):

$$E(\mathbf{w}, w_0 | \mathcal{X}) = - \sum_t r^t \log y^t + (1 - r^t) \log(1 - y^t) \tag{10.23}$$

我们使用梯度下降法最小化互熵, 等价于最大化似然或对数似然。如果 $y = \text{sigmoid}(a) = 1/(1 + \exp(-a))$, 则它的导数为

$$\frac{dy}{da} = y(1 - y)$$

并且我们得到如下更新方程:

$$\begin{aligned}
 \Delta \mathbf{w}_j &= -\eta \frac{\partial E}{\partial w_j} = \eta \sum_t \left(\frac{r^t}{y^t} - \frac{1 - r^t}{1 - y^t} \right) y^t (1 - y^t) \mathbf{x}_j^t \\
 &= \eta \sum_t (r^t - y^t) \mathbf{x}_j^t, j = 1, \dots, d \\
 \Delta w_0 &= -\eta \frac{\partial E}{\partial w_0} = \eta \sum_t (r^t - y^t)
 \end{aligned} \tag{10.24}$$

最好用接近于 0 的随机值初始化 w_j ; 通常, 它们从区间 $[-0.01, 0.01]$ 中均匀地抽取。这样做的理由是, 如果 w_j 数值很大, 则加权和平可能也很大并且 S 形函数可能饱和。从图 10-5 我们看到, 如果初始权重接近于 0, 则和在区域中间, 那里导数非零, 可以进行更新。如果加权和平很大 (小于 -5 或大于 +5), 则 S 形函数的导数几乎为 0, 权值将不会更新。

伪代码在图 10-6 中。我们看图 10-7 中的例子，其中输入是一维的。直线 $w x + w_0$ 和它的 S 形函数之后的值都作为学习迭代次数的函数显示。我们看到，为了得到输出 0 和 1，S 形函数逐渐适应，这通过增大 w 的数值实现。

```

For  $j = 0, \dots, d$ 
   $w_j \leftarrow \text{rand}(-0.01, 0.01)$ 
Repeat
  For  $j = 0, \dots, d$ 
     $\Delta w_j \leftarrow 0$ 
  For  $t = 1, \dots, N$ 
     $o \leftarrow 0$ 
    For  $j = 0, \dots, d$ 
       $o \leftarrow o + w_j x_j^t$ 
     $y \leftarrow \text{sigmoid}(o)$ 
    For  $j = 0, \dots, d$ 
       $\Delta w_j \leftarrow \Delta w_j + (r^t - y) x_j^t$ 
    For  $j = 0, \dots, d$ 
       $w_j \leftarrow w_j + \eta \Delta w_j$ 
Until 收敛
  
```

图 10-6 对于具有两个类、单个输出，实现梯度下降的逻辑斯谛判别分析算法。对于 w_0 ，我们假定存在一个附加的输入 x_0 ，它总为 +1: $x_0^t \equiv +1, \forall t$

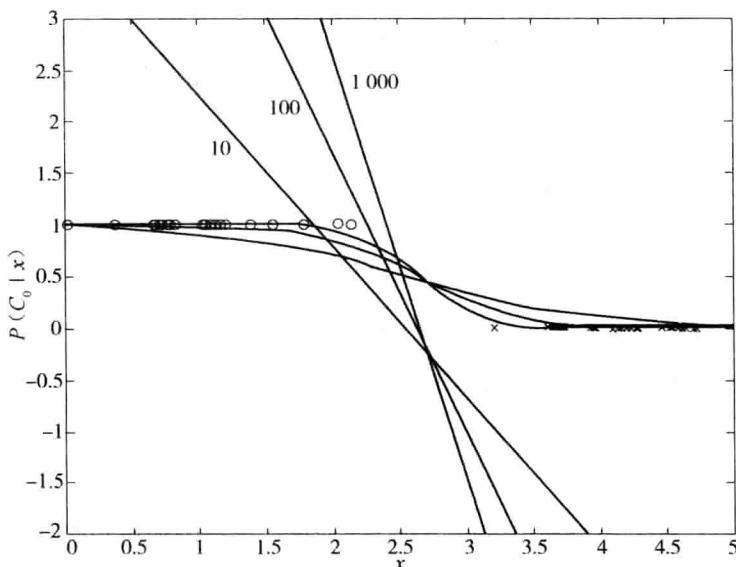


图 10-7 对于一元两类问题(用“ \circ ”和“ \times ”显示)，样本上 10 次、100 次和 1000 次迭代之后，直线 $w x + w_0$ 和 S 形函数输出的演变

一旦训练完成并且我们得到了最终的 w 和 w_0 ，在检验阶段，给定 x ，我们计算 $y = \text{sigmoid}(w^T x + w_0)$ ，并且如果 $y > 0.5$ 则选择 C_1 ，否则选择 C_2 。这意味，为了最小化误分类数，我们不需要学习到 y 是 0 或 1，而只需要学习到 y 小于或大于 0.5。如果超过该点后我们还继续学习，互熵将继续降低 ($|w_j|$ 将继续增加，硬化 S 形函数)，但是误分类数将不会

减少。通常，我们一直训练，直到误分类数不再减少（如果类是线性可分的，它将为 0）。实际上，在达到零训练误差之前提前停止是一种正则化形式。因为以权重几乎为 0 开始，并且它们随着训练继续而远离 0，所以提前停止对应于具有更多接近于零权重而实际上参数更少的模型。

注意，尽管为了导出判别式，我们假定类密度的对数比是线性的，但是我们直接估计后验概率，而不显式地估计 $p(\mathbf{x} | C_i)$ 或 $P(C_i)$ 。

10.7.2 多类问题

现在，让我们推广到 $K > 2$ 个类：我们取其中一个类，例如 C_k ，作为引述类并假定

$$\log \frac{p(\mathbf{x} | C_i)}{p(\mathbf{x} | C_K)} = \mathbf{w}_i^T \mathbf{x} + w_{i0} \quad (10.25)$$

于是，我们有

$$\frac{P(C_i | \mathbf{x})}{P(C_K | \mathbf{x})} = \exp[\mathbf{w}_i^T \mathbf{x} + w_{i0}] \quad (10.26)$$

其中 $w_{i0} = w_{i0}^0 + \log P(C_i)/P(C_K)$ 。

我们看到

$$\begin{aligned} \sum_{i=1}^{K-1} \frac{P(C_i | \mathbf{x})}{P(C_K | \mathbf{x})} &= \frac{1 - P(C_K | \mathbf{x})}{P(C_K | \mathbf{x})} = \sum_{i=1}^{K-1} \exp[\mathbf{w}_i^T \mathbf{x} + w_{i0}] \\ \Rightarrow P(C_K | \mathbf{x}) &= \frac{1}{1 + \sum_{i=1}^{K-1} \exp[\mathbf{w}_i^T \mathbf{x} + w_{i0}]} \end{aligned} \quad (10.27)$$

并且还有

$$\begin{aligned} \frac{P(C_i | \mathbf{x})}{P(C_K | \mathbf{x})} &= \exp[\mathbf{w}_i^T \mathbf{x} + w_{i0}] \\ \Rightarrow P(C_i | \mathbf{x}) &= \frac{\exp[\mathbf{w}_i^T \mathbf{x} + w_{i0}]}{1 + \sum_{j=1}^{K-1} \exp[\mathbf{w}_j^T \mathbf{x} + w_{j0}]}, i = 1, \dots, K-1 \end{aligned} \quad (10.28)$$

为了一致地处理所有的类，我们记作

$$y_i = \hat{p}(C_i | \mathbf{x}) = \frac{\exp[\mathbf{w}_i^T \mathbf{x} + w_{i0}]}{\sum_{j=1}^K \exp[\mathbf{w}_j^T \mathbf{x} + w_{j0}]}, i = 1, \dots, K \quad (10.29)$$

这称为软最大 (softmax) 函数 (Bridle 1990)。如果一个类的加权和明显大于其他类的加权和，则通过取指数和规范化推进之后，它对应的 y_i 将接近于 1，而其他的将接近于 0。这样，除了可导之外，它就像取最大，因此得名软最大。软最大还保证 $\sum_i y_i = 1$ 。

让我们看如何学习参数：在 $K > 2$ 个类这种情况下，每个样本点是一次多项试验取值，即 $\mathbf{r}^t | \mathbf{x}^t \sim \text{Mult}_K(1, \mathbf{y}^t)$ ，其中 $y_i^t \equiv P(C_i | \mathbf{x}^t)$ 。样本的似然为

$$l(\{\mathbf{w}_i, w_{i0}\}_i | \mathcal{X}) = \prod_i \prod_i (y_i^t)^{r_i^t} \quad (10.30)$$

而误差函数又是互熵：

222
223

224

$$E(\{\mathbf{w}_i, w_{j0}\}_i | \mathcal{X}) = - \sum_t \sum_i r_i^t \log y_i^t \quad (10.31)$$

我们再次使用梯度下降。如果 $y_i = \exp(a_i) / \sum_j \exp(a_j)$ ，则我们有

$$\frac{\partial y_i}{\partial a_j} = y_i (\delta_{ij} - y_j) \quad (10.32)$$

其中 δ_{ij} 是 Kronecker δ ，如果 $i=j$ 它为 1，如果 $i \neq j$ 它为 0 (习题 3)。给定 $\sum_i r_i^t = 1$ ，对于 $j = 1, \dots, K$ ，我们有如下更新方程

$$\begin{aligned} \Delta \mathbf{w}_j &= \eta \sum_t \sum_i \frac{r_i^t}{y_i^t} (\delta_{ij} - y_j^t) \mathbf{x}^t \\ &= \eta \sum_t \sum_i r_i^t (\delta_{ij} - y_j^t) \mathbf{x}^t \\ &= \eta \sum_t \left[\sum_i r_i^t \delta_{ij} - y_j^t \sum_i r_i^t \right] \mathbf{x}^t \\ &= \eta \sum_t (r_j^t - y_j^t) \mathbf{x}^t \\ \Delta w_{j0} &= \eta \sum_t (r_j^t - y_j^t) \end{aligned} \quad (10.33)$$

注意，由于软最大中的规范化， w_j 和 w_{j0} 不仅受 $\mathbf{x}^t \in C_j$ 的影响，而且还受 $\mathbf{x}^t \in C_i (i \neq j)$ 的影响。更新判别式使得取软最大之后正确的类具有最大的加权和，而其他类的加权和尽可能小。伪代码在图 10-8 中给出。对于具有三个类的二维样本，轮廓线在图 10-9 中给出，而判别式和后验概率在图 10-10 中。

```

For  $i=1, \dots, K$ , For  $j=0, \dots, d$ ,  $w_{ij} \leftarrow \text{rand}(-0.01, 0.01)$ 
Repeat
  For  $i=1, \dots, K$ , For  $j=0, \dots, d$ ,  $\Delta w_{ij} \leftarrow 0$ 
  For  $t=1, \dots, N$ 
    For  $i=1, \dots, K$ 
       $o_i \leftarrow 0$ 
      For  $j=0, \dots, d$ 
         $o_i \leftarrow o_i + w_{ij} x_j^t$ 
      For  $i=1, \dots, K$ 
         $y_i \leftarrow \exp(o_i) / \sum_k \exp(o_k)$ 
      For  $i=1, \dots, K$ 
        For  $j=0, \dots, d$ 
           $\Delta w_{ij} \leftarrow \Delta w_{ij} + (r_i^t - y_i^t) x_j^t$ 
      For  $i=1, \dots, K$ 
        For  $j=0, \dots, d$ 
           $w_{ij} \leftarrow w_{ij} + \eta \Delta w_{ij}$ 
  Until 收敛

```

图 10-8 对于 $K > 2$ 个类，实现梯度下降的逻辑斯谛判别式算法。
为了一般起见，对于任意 t ，我们取 $x_0^t = 1$

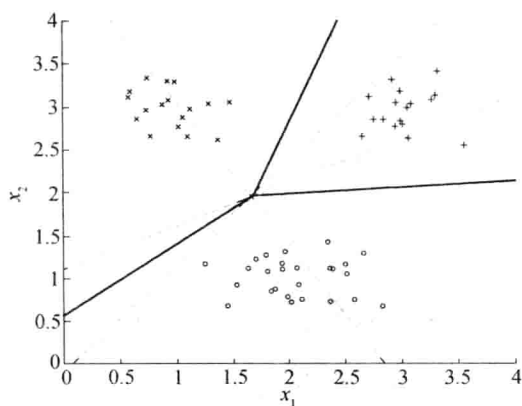


图 10-9 对于具有三个类的二维问题，逻辑斯谛判别式发现的解。细线是 $g_i(\mathbf{x}) = 0$ ，而粗线是取极大的线性分类器得到的边界

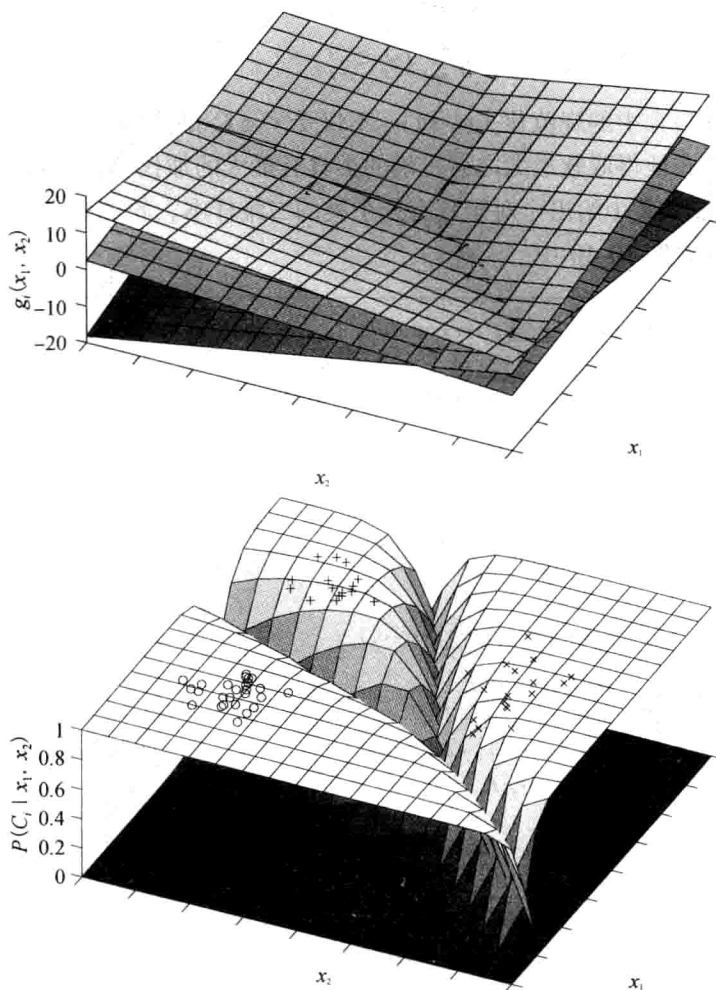


图 10-10 对于图 10-9 中的相同样本，线性判别式(上)和软最大后的后验概率(下)

在检验阶段, 我们计算所有的 $y_k, k=1, \dots, K$, 并且如果 $y_i = \max_k y_k$, 则选择 C_i 。我们仍然不必为尽可能地极小化互熵而一直训练; 我们只需要训练到正确的类具有最大的加权和, 并通过检查误分类数, 提前停止训练。

当数据是正态分布时, 逻辑斯谛判别式与参数的、基于正态的线性判别式具有大致相当的错误率 (McLachlan 1992)。当类条件密度非正态时, 或当它们非单峰时, 只要类是线性可分的, 逻辑斯谛判别式仍然可以使用。

当然, 类条件密度的比不限于是线性的 (Anderson 1982; McLachlan 1992)。假定一个二次判别式, 我们有

$$\log \frac{p(\mathbf{x} | C_i)}{p(\mathbf{x} | C_k)} = \mathbf{x}^T \mathbf{W}_i \mathbf{x} + \mathbf{w}_i^T \mathbf{x} + w_{i0} \quad (10.34)$$

对应并推广具有不同协方差矩阵、满足多元正态类条件分布的参数判别式。当 d 很大时, 正如我们可以化简 (正规化) Σ_i 一样, 我们可以通过只考虑它的前面的本征向量, 对 \mathbf{W}_i 做同样的事。

正如 10.2 节所讨论的, 可以用基本变量的任意指定函数作为 x -变量。例如, 我们可以把判别式写成非线性基函数的线性和

$$\log \frac{p(\mathbf{x} | C_i)}{p(\mathbf{x} | C_k)} = \mathbf{w}_i^T \boldsymbol{\phi}(\mathbf{x}) + w_{i0} \quad (10.35)$$

其中, $\boldsymbol{\phi}(\cdot)$ 是基函数, 可以看作变换后的变量。用神经网络的术语, 这称作多层感知器 (multilayer perception) (第 11 章), 而 S 形函数是最常用的基函数。当使用高斯基函数时, 这种模型称作径向基函数 (radial basis function) (第 12 章)。我们甚至可以使用完全非参数的方法, 如 Parzen 窗口 (第 8 章)。

10.8 回归判别式

对于回归, 概率模型是

$$r^i = y^i + \varepsilon \quad (10.36)$$

其中 $\varepsilon \sim \mathcal{N}(0, \sigma^2)$ 。如果 $r^i \in \{0, 1\}$, 使用 S 形函数, 则 y^i 可能限于落在该区间。假定线性模型和两个类, 我们有

$$y^i = \text{sigmoid}(\mathbf{w}^T \mathbf{x}^i + w_0) = \frac{1}{1 + \exp[-(\mathbf{w}^T \mathbf{x}^i + w_0)]} \quad (10.37)$$

假定 $r | \mathbf{x} \sim \mathcal{N}(y, \sigma^2)$, 则回归中的样本似然为

$$l(\mathbf{w}, w_0 | \mathcal{X}) = \prod_i \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{(r^i - y^i)^2}{2\sigma^2}\right] \quad (10.38)$$

最大化该对数似然是最小化误差的平方和:

$$E(\mathbf{w}, w_0 | \mathcal{X}) = \frac{1}{2} \sum_i (r^i - y^i)^2 \quad (10.39)$$

使用梯度下降, 我们得到

$$\begin{aligned} \Delta \mathbf{w} &= \eta \sum_i (r^i - y^i) y^i (1 - y^i) \mathbf{x}^i \\ \Delta w_0 &= \eta \sum_i (r^i - y^i) y^i (1 - y^i) \end{aligned} \quad (10.40)$$

当存在 $K > 2$ 个类时,也可以使用这种方法。概率模型是

$$\mathbf{r}^t = \mathbf{y}^t + \boldsymbol{\varepsilon} \quad (10.41)$$

其中 $\boldsymbol{\varepsilon} \sim \mathcal{N}_k(0, \sigma^2 \mathbf{I}_K)$ 。假定每个类一个线性模型,我们有

$$y_i^t = \text{sigmoid}(w_i^T \mathbf{x}^t + w_{i0}) = \frac{1}{1 + \exp[-(w_i^T \mathbf{x}^t + w_{i0})]} \quad (10.42)$$

于是,样本的似然为

$$l(\{\mathbf{w}_i, w_{i0}\}_i | \chi) = \prod_i \frac{1}{(2\pi)^{K/2} |\Sigma|^{1/2}} \exp\left[-\frac{\|\mathbf{r}^t - \mathbf{y}^t\|^2}{2\sigma^2}\right] \quad (10.43)$$

而误差函数为

$$E(\{\mathbf{w}_i, w_{i0}\}_i | \chi) = \frac{1}{2} \sum_i \|\mathbf{r}^t - \mathbf{y}^t\|^2 = \frac{1}{2} \sum_i \sum_i (r_i^t - y_i^t)^2 \quad (10.44)$$

对于 $i = 1, \dots, K$, 更新方程为

$$\begin{aligned} \Delta \mathbf{w}_i &= \eta \sum_i (r_i^t - y_i^t) y_i^t (1 - y_i^t) \mathbf{x}^t \\ \Delta w_{i0} &= \eta \sum_i (r_i^t - y_i^t) y_i^t (1 - y_i^t) \end{aligned} \quad (10.45)$$

注意,在这样做时,我们并未使用 y_i 中的一个为 1,其余为 0,或 $\sum_i y_i = 1$ 的信息。由于输出是类后验概率估计,式(10.29)的软最大函数使得我们可以纳入这些附加信息。在 $K > 2$ 的情况下,使用 S 形函数,我们像对待独立函数一样对待 y_i 。

还要注意,对于给定的类,如果我们使用回归方法,则更新将进行到正确的输出为 1,其余为 0 才停止。事实上,这是不必要的,因为在检验时,我们只是选择最大的。训练到正确的输出大于其他输出就足够了,这正是软最大函数所做的。

因此,当类不是互斥的和穷举的时,也就是说,对于一个 \mathbf{x}^t ,所有的 r_i^t 可能都为 0,即 \mathbf{x}^t 不属于任何一个类;或者当类重叠时,可能多个 r_i^t 为 1 时,这种具有多个 S 形函数的方法更可取。

10.9 注释

由于其简单性,线性判别分析是模式识别研究最多的课题(Duda、Hart 和 Stork 2001; McLachlan 1992)。我们在第 4 章讨论了具有公共协方差矩阵的高斯分布情况,在第 6 章讨论了费希尔线性判别式,并在本章讨论了逻辑斯谛判别式。在第 11 章,我们将讨论感知器,它是线性判别式的神经网络实现。第 13 章我们将讨论支持向量机,这是另一种类型的线性判别式。

逻辑斯谛判别式的更详细讨论在 Anderson 1982 和 McLachlan 1992 中。逻辑斯谛(S 形)函数是分对数的逆,在伯努利抽样中称作规范链(canonical link)。软最大是对多元正态抽样的拓广。关于广义线性模型(generalized linear model)的更多信息在 McCulloch 和 Nelder 1989 中。

使用非线性基函数推广线性模型是一种历史悠久的想法。我们将讨论多层感知器(参见第 11 章)和径向基函数(参见第 12 章),其中基函数的参数也可以在学习判别式时由数据学习。支持向量机(参见第 13 章)使用由这些基函数形成的核函数。

10.10 习题

1. 对于如下每个基函数，指出它在何处非零：

a. $\sin(x_1)$	b. $\exp(-(x_1 - a)^2/c)$	c. $\exp(-\ \mathbf{x} - \mathbf{a}\ ^2/c)$
d. $\log(x_2)$	e. $1(x_1 > c)$	f. $1(ax_1 + bx_2 > c)$
2. 对于图 10-2 的二维情况，证明式 (10.4) 和式 (10.5)。
3. 证明软最大 $y_i = \exp(a_i) / \sum_j \exp(a_j)$ 的导函数是 $\partial y_i / \partial a_i = y_i(\delta_{ij} - y_j)$ ，其中如果 $i = j$ 则 δ_{ij} 为 1，否则 δ_{ij} 为 0。
4. 令 $K = 2$ ，证明用两个软最大输出等于使用一个 S 形输出。
5. 在式 (10.34) 中，我们如何学习 \mathbf{W}_i ？
6. 当像式 (10.34) 中那样使用二次 (或更高阶) 判别式时，如何保持方差受控？
7. 在梯度下降时，对所有的 x_j 使用单个 η 意味着什么？
8. 在单变量情况下，对于如图 10.7 中的分类， w 和 w_0 对应于什么？
9. 假设对于单变量 x ， $x \in (2, 4)$ 属于 C_1 ，而 $x < 2$ 或 $x > 4$ 属于 C_2 。如何使用线性判别式把这两个类分离？

230

10.11 参考文献

- Aizerman, M. A., E. M. Braverman, and L. I. Rozonoer. 1964. "Theoretical Foundations of the Potential Function Method in Pattern Recognition Learning." *Automation and Remote Control* 25: 821-837.
- Anderson, J. A. 1982. "Logistic Discrimination." In *Handbook of Statistics*, Vol. 2, *Classification, Pattern Recognition and Reduction of Dimensionality*, ed. P. R. Krishnaiah and L. N. Kanal, 169-191. Amsterdam: North Holland.
- Bridle, J. S. 1990. "Probabilistic Interpretation of Feedforward Classification Network Outputs with Relationships to Statistical Pattern Recognition." In *Neurocomputing: Algorithms, Architectures and Applications*, ed. F. Fogelman-Soulie and J. Herault, 227-236. Berlin: Springer.
- Duda, R. O., P. E. Hart, and D. G. Stork. 2001. *Pattern Classification*, 2nd ed. New York: Wiley.
- McCullagh, P., and J. A. Nelder. 1989. *Generalized Linear Models*. London: Chapman and Hall.
- McLachlan, G. J. 1992. *Discriminant Analysis and Statistical Pattern Recognition*. New York: Wiley.
- Vapnik, V. 1995. *The Nature of Statistical Learning Theory*. New York: Springer.

231

多层感知器

多层感知器是一种人工神经网络结构，是非参数估计器，可以用于分类和回归。我们讨论为各种应用训练多层感知器的后向传播算法。

11.1 引言

人工神经网络模型，其中之一是我们本章讨论的多层感知器，其灵感源于模拟人脑。旨在理解人脑功能，并朝着这一目标努力的认知科学家和神经学家 (Posner 1989) 构建了人脑的神经网络模型，并开展了模拟研究。

然而，在工程上，我们的目标不是理解人脑的本质，而是构建有用的机器。我们对人工神经网络 (artificial neural networks) 感兴趣，因为我们相信它们可能帮助我们建立更好的计算机系统。人脑是一种信息处理装置，具有非凡的能力并且在许多领域，例如，视觉、语音识别和学习方面，都超过了当前的工程产品。如果在机器上实现，这些应用显然都具有经济效益。如果我们能够理解人脑如何实现这些功能，我们就可以用形式算法定义这些任务的解，并且在计算机上实现它们。

人脑与计算机很不相同。计算机通常只有一个处理器，而人脑却包含大量 (10^{11} 个) 并行操作的处理单元，称作神经元 (neuron)。尽管处理细节尚不清楚，但是人们相信这些处理单元比计算机中的处理器简单得多，并且慢得多。使得人脑不同寻常并且被认为提供了其计算能力的是连通性：人脑的神经元具有连接，称作突触 (synapse)，连接到大约 10^4 个其他神经元，所有神经元都并行地操作。在计算机中，处理器是主动的，而存储是分散和被动的，但是我们认为在人脑中，处理和存储都在网络上分布；处理由神经元来做，而记忆在神经元之间的突触中。

11.1.1 理解人脑

根据 Marr (1982)，理解一个信息处理系统具有三个层面，称作分析层面 (levels of analysis)：

1) 计算理论 (computational theory) 对应计算目标和任务的抽象定义。

2) 表示和算法 (representation and algorithm) 是关于输入和输出如何表示和从输入到输出变换的算法说明。

3) 硬件实现 (hardware implementation) 是系统的实际物理实现。

一个例子是排序：计算理论是对给定的元素集合排序。表示可以使用整数，而算法可以是 Quicksort (快速排序)。编译后，特定处理机对二进制表示的整数排序的可执行代码是一种硬件实现。

基本思想是，对于相同的计算理论，可以有多种表示和在相应表示上操控符号的算法。

类似地，对于给定的表示和算法，可以有多种硬件实现。我们可以使用众多排序算法中的一种，并且即使相同的算法也可以在使用不同处理器的计算机上编译，导致不同的硬件实现。

考虑另一个例子，“6”、“VI”和“110”是数字 6 的不同表示。加法的不同算法依赖于所使用的表示。数字计算机使用二进制表示，并具有这种表示的加法电路，这是一种特殊的硬件实现。在算盘上，数用不同的方法表示，并且加法对应不同的指令集，这是另一种硬件实现。当我们在大脑中将两个数相加时，我们使用另一种表示和一种适合于这种表示的算法，这由神经元实现。但是，所有不同的硬件实现（例如，我们、算盘和数字计算机）都实现了相同的计算理论——加法。

经典的例子是自然和人工的飞行器之间的不同：麻雀拍打它的双翼；商用飞机并不拍打机翼，而是使用喷气引擎。麻雀和飞机是两种硬件实现，为不同的目的而构建，满足不同的约束。但是它们都实现了相同的理论——空气动力学。

人脑是学习或模式识别的一种硬件实现。如果从这种特定的实现，我们可以做逆工程，提取人脑使用的表示和算法，并且如果我们能够从中获得计算理论，则我们可以使用另一种表示和算法，然后得到更适合我们的含义和约束的硬件实现。我们希望我们的实现价格低廉、快速，并且更准确。

当初构建飞行器时，直到发现空气动力学之前，我们一直在尝试构建看上去非常像鸟的飞行器。与此相同，直到我们发现智能的计算理论之前，早期尝试构建具有大脑能力的结构看上去将很像大脑，是具有大量处理单元的网络。因此可以说，就理解大脑而言，当我们研究人工神经网络时，我们处于表示和算法层面。

正如羽毛与飞行不相关一样，将来我们可能发现神经元和突触与智能并无关系。但是，在此之前，我们对理解大脑机能感兴趣还有另一个原因，这种原因与并行处理有关。

11.1.2 神经网络作为并行处理的典范

自 20 世纪 80 年代以来，具有数以千计处理器的计算机系统已经商品化。然而，用于这种并行结构的软件并不像硬件发展这么快。原因是到目前为止我们的计算理论几乎都基于串行的、单处理器机器。我们不能有效地使用并行机，因为我们不能有效地对它们编程。

234
235

主要有两种并行处理(parallel processing)范型：在单指令多数据(SIMD)机，所有的处理器都执行相同的指令，但是在不同的数据上执行。在多指令多数据(MIMD)机，不同的处理器可以在不同的数据上执行不同的指令。SIMD 机容易编程，因为只需要写一个程序。然而，问题很少具有这种有规律的结构以能在 SIMD 机上并行地执行。MIMD 机更一般，但是为每个处理器编写单独的程序并不是一件容易的任务；其他问题涉及同步、处理器之间的数据传输等。SIMD 机也比较容易构建，并且如果它们都是 SIMD 机，则可以构建具有更多处理器的机器。在 MIMD 机中，处理器更加复杂，并且还要为处理器任意地交换数据构建更复杂的通信网络。

现在，假设我们可以有机器，其中处理器比 SIMD 处理器复杂一点，但没有 MIMD 处理器复杂。假定我们有一些简单处理器，具有少量局部存储器，可以存放一些参数。每个处理器实现一个固定的函数，并且执行与 SIMD 处理器一样的指令；但是通过将不同的值装入局部存储器，它们可以做不同的事情，并且整个操作可以在这些处理器上分布执行。这样，我

们将有可以称作神经指令多数据(NIMD)机,其中每个处理器对应一个神经元,局部参数对应它的突触权重,而整个结构是一个神经网络。如果每个处理器中实现的功能很简单,并且局部存储器很小,则许多这样的处理器可以放在一个芯片中。

现在的问题是将任务分布到这种处理器的网络中和确定局部参数的值。这是学习进行的地方:如果这样的机器可以从实例学习,则我们自己不需要为这种机器编制程序和决定参数值。

因此,人工神经网络是一种我们可以使用当前技术构建的、利用并行硬件的方法——多亏了学习——它们不需要编程。因此,我们也不必费神为它们编程。

本章,我们讨论这种结构和如何训练它们。记住,人工神经网络操作是一种数学函数,它们可以在串行计算机上实现,并且训练网络与我们在前面的章节中讨论的统计学技术并无太大差别。仅当我们有并行硬件,并且仅当网络太大,不能在串行机上快速模拟时,考虑这些操作在简单处理单元的网络上进行才是有意义的。

11.2 感知器

感知器(perceptron)是基本处理元素。它具有输入,其输入可能来自环境或者可以是其他感知器的输出。与每个输入 $x_j \in \mathfrak{R} (j=1, \dots, d)$ 相关联的是一个连接权重(connection weight)或突触权重(synaptic weight) $w_j \in \mathfrak{R}$, 而输出 y 在最简单情况下是输入的加权和(参见图 11-1):

$$y = \sum_{j=1}^d w_j x_j + w_0 \quad (11.1)$$

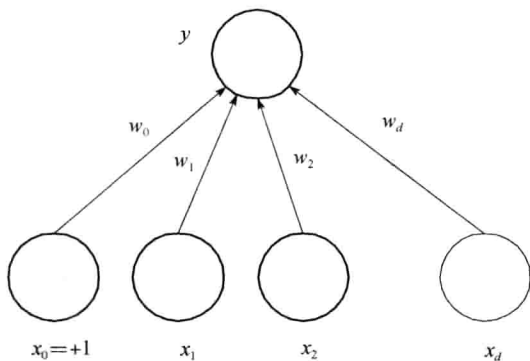


图 11-1 简单感知器, $x_j (j=1, \dots, d)$ 是输入单元, x_0 是偏倚单元, 其值总是为 1。 y 是输出单元。 w_j 是从 x_j 到输出的有向连接的权重

其中 w_0 是截距值, 它使模型更一般; 通常把它作为一个来自附加的偏倚单元(bias unit) x_0 的权重, 而 x_0 总是为 +1。我们可以把感知器的输出写成点积

$$y = \mathbf{w}^T \mathbf{x} \quad (11.2)$$

其中 $\mathbf{w} = [w_0, w_1, \dots, w_d]^T$, $\mathbf{x} = [1, x_1, \dots, x_d]^T$ 是增广向量(augmented vector), 包含偏倚权重和输入。

在检验时, 给定权重 \mathbf{w} , 对于输入 \mathbf{x} 我们计算输出 y 。为了实现给定的任务, 我们需要学习系统的参数权重 \mathbf{w} , 使得我们可以产生给定输入的正确输出。

当 $d=1$ 并且 x 通过输入单元由环境馈入时, 我们有

$$y = wx + w_0$$

这是以 w 为斜率、 w_0 为截距的直线方程。这样，这种具有一个输入和一个输出的感知器可以用来实现线性拟合。使用多个输入，直线变成了(超)平面，而具有多个输入的感知器可以实现多元线性拟合。给定样本，通过回归可以找出参数 w_j (见 5.8 节)。

(11.1) 式定义的感知器定义了一个超平面，因此可以用来将输入空间划分成两部分： y 值为正的半个空间和 y 值为负的半个空间(参见第 10 章)。通过用它实现线性判别函数，检查输出的符号，感知器可以将两个类分开。如果我们定义 $s(\cdot)$ 为阈值函数(threshold function)

$$s(a) = \begin{cases} 1 & \text{如果 } a > 0 \\ 0 & \text{否则} \end{cases} \quad (11.3)$$

则如果 $s(\mathbf{w}^T \mathbf{x}) > 0$ 我们可以选择 C_1 ，否则选择 C_2 。

记住使用线性判别式假定类是线性可分的。也就是说，假定可以找到分开 $\mathbf{x}' \in C_1$ 和 $\mathbf{x}' \in C_2$ 的超平面 $\mathbf{w}^T \mathbf{x} = 0$ 。如果在后一阶段我们需要后验概率(例如计算风险)，我们需要在输出上使用 S 型函数

$$o = \mathbf{w}^T \mathbf{x} \\ y = \text{sigmoid}(o) = \frac{1}{1 + \exp[-\mathbf{w}^T \mathbf{x}]} \quad (11.4) \quad \boxed{238}$$

当存在 $K > 2$ 个输出时，有 K 个感知器，每个都具有权重向量 \mathbf{w}_i (参见图 11-2)

$$y_i = \sum_{j=1}^d w_{ij} x_j + w_{i0} = \mathbf{w}_i^T \mathbf{x} \\ \mathbf{y} = \mathbf{W} \mathbf{x} \quad (11.5)$$

其中 w_{ij} 是从输入 x_j 到输出 y_i 的连接权重。 \mathbf{W} 是 w_{ij} 的 $K \times (d+1)$ 矩阵，其行是 K 个感知器的权重向量。当用于分类时，在检验阶段，如果 $y_i = \max_k y_k$ ，则我们选择 C_i 。

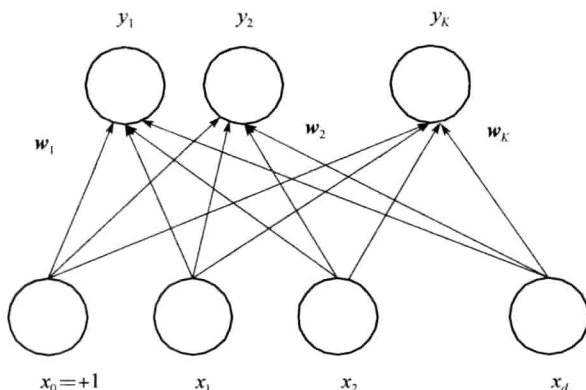


图 11-2 K 个并行的感知器。 $x_j (j=0, \dots, d)$ 是输入， $y_i (i=1, \dots, K)$ 是输出。 w_{ij}^{\ominus}

是从输入 x_j 到输出 y_i 的连接权重。每个输出都是输入的加权和。当用于 K

类问题时，有一个后处理，选择最大的，或者需要后验概率时选择软最大
在神经网络中，每个感知器的值是它的输入和它的突触权重的局部函数(local function)。

\ominus 图中并未标出 w_{ij} ，但标出了权重向量 $\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_K$ 。——取自作者勘误

然而在分类中, 如果我们需要后验概率(而不仅是获胜者类的编码)并使用软最大, 则我们还需要其他输出值。这样, 为了将其作为神经网络实现, 我们可以将其看作一个两阶段过程, 其中

[239]
$$o_i = \mathbf{w}_i^T \mathbf{x}$$

$$y_i = \frac{\exp o_i}{\sum_k \exp o_k} \quad (11.6)$$

回忆一下, 通过定义附加的输入, 例如, 定义 $x_3 = x_1^2$, $x_4 = x_2^2$, $x_5 = x_1 x_2$ (10.2 节), 线性模型也可以用于多项式逼近。对于感知器也可以这样做 (Durbin 和 Rumelhart 1989)。在 11.5 节, 我们将看到多层感知器, 那里非线性函数从数据中学习, 而不是先验假定。

第 10 章讨论的线性判别式的任何方法都可以离线地计算 \mathbf{w}_i , $i = 1, \dots, K$, 然后插入到网络中。这包括具有公共协方差矩阵的参数方法、逻辑斯谛判别式、借助于回归的判别式和支持向量机。在某些情况下, 在训练开始时我们并没有全部样本, 并且随着新的实例到来, 我们需要迭代地更新参数; 我们将在 11.3 节讨论这种在线学习。

式(11.5)定义了一个从 d -维空间到 K -维空间的变换, 如果 $K < d$, 它也可以用于维归约。我们可以使用第 6 章中的任何方法(如 PCA)离线地计算 \mathbf{W} , 然后使用感知器实现变换。在这种情况下, 我们有两层网络, 其中第一层感知器实现线性变换, 而第二层在新空间实现线性回归或分类。注意, 由于两层都是线性变换, 因此它们可以组合并用一层表示。在 11.5 节, 我们将看到更有趣的情况, 其中第一层实现非线性维归约。

11.3 训练感知器

感知器定义了一个超平面, 而神经网络感知器只不过是实现超平面的一种方法。给定数据样本, 权重可以离线地计算, 并且将它们代入后, 感知器就可以用来计算输出的值。

[240] 在训练神经网络时, 如果未提供全部样本而是逐个提供实例, 则我们通常使用在线学习, 并且在每个实例之后更新网络参数, 让网络缓慢地及时调整。这种方法是令人感兴趣的, 有如下原因:

1) 这使得我们不必在外存存放训练样本, 不必在优化时存放中间结果。对于大样本, 像支持向量机这样的方法(10.9 节)开销可能相当高, 而对于某些应用, 我们可能更愿意选择较简单的方法, 不必存放全部样本并在其上求解复杂的优化问题。

2) 问题可能随时变化, 这意味样本的分布不固定, 训练集不能预先选定。例如, 我们可能正在实现一个语音识别系统, 它要自动适应它的用户。

3) 可能存在系统的物理变化。例如, 在一个机器人系统中, 系统部件可能磨损, 传感器可能失灵。

对于在线学习(online learning), 我们不需要全部样本而是需要单个实例上的误差函数。从随机初始权重开始, 在每次迭代中, 我们都对参数稍加调整, 以最小化误差, 而不忘记我们先前学到的。如果误差函数是可微的, 则我们可以使用梯度下降。

例如, 对于回归, 单个具有标引 t 的实例 (\mathbf{x}^t, r^t) 的误差为

$$E^t(\mathbf{w} | \mathbf{x}^t, r^t) = \frac{1}{2}(r^t - y^t)^2 = \frac{1}{2}[r^t - (\mathbf{w}^T \mathbf{x}^t)]^2$$

并且对于 $j=0, \dots, d$, 在线更新为

$$\Delta w_j^t = \eta(r^t - y^t)x_j^t \quad (11.7)$$

其中 η 是学习因子, 随收敛逐渐减小。这称作随机梯度下降(stochastic gradient descent)。

类似地, 可以对使用逻辑斯谛判别式的分类问题导出更新规则。在那里, 每个模式之后进行更新, 而不是把它们累加在一起, 在完全扫描整个训练集之后再进行更新。对于两个类, 对单个实例 $(\mathbf{x}^t, \mathbf{r}^t)$, 其中如果 $\mathbf{x}^t \in C_1$ 则 $r_i^t = 1$, 如果 $\mathbf{x}^t \in C_2$ 则 $r_i^t = 0$, 单个输出为

$$y^t = \text{sigmoid}(\mathbf{w}^T \mathbf{x}^t)$$

241

而互熵为

$$E^t(\{\mathbf{w}_i\}_i | \mathbf{x}^t, \mathbf{r}^t) = -\mathbf{r}^t \log y^t + (1 - \mathbf{r}^t) \log(1 - y^t)$$

使用梯度下降, 对于 $j=0, \dots, d$, 我们得到如下更新规则:

$$\Delta w_j^t = \eta(r^t - y^t)x_j^t \quad (11.8)$$

当存在 $K > 2$ 个类时, 对单个实例 $(\mathbf{x}^t, \mathbf{r}^t)$, 其中如果 $\mathbf{x}^t \in C_i$ 则 $r_i^t = 1$, 否则 $r_i^t = 0$, 输出为

$$y_i^t = \frac{\exp \mathbf{w}_i^T \mathbf{x}^t}{\sum_k \exp \mathbf{w}_k^T \mathbf{x}^t}$$

而互熵为

$$E^t(\{\mathbf{w}_i\}_i | \mathbf{x}^t, \mathbf{r}^t) = -\sum_i r_i^t \log y_i^t$$

使用梯度下降, 对于 $i=1, \dots, K, j=0, \dots, d$, 我们得到如下更新规则:

$$\Delta w_{ij}^t = \eta(r_i^t - y_i^t)x_j^t \quad (11.9)$$

除了不在所有的实例上求和, 而是在单个实例之后更新外, 这与我们在 10.7 节中看到的方程一样。算法的伪代码在图 11-3 中, 它是图 10-8 算法的在线版本。

```

For  $i = 1, \dots, K$ 
  For  $j = 0, \dots, d$ 
     $w_{ij} \leftarrow \text{rand}(-0.01, 0.01)$ 
  Repeat
    For 随机序下的所有  $(\mathbf{x}^t, \mathbf{r}^t) \in \mathcal{X}$ 
      For  $i = 1, \dots, K$ 
         $o_i \leftarrow 0$ 
        For  $j = 0, \dots, d$ 
           $o_i \leftarrow o_i + w_{ij}x_j^t$ 
        For  $k = 1, \dots, K$ 
           $y_k \leftarrow \exp(o_k) / \sum_k \exp(o_k)$ 
        For  $i = 1, \dots, K$ 
          For  $j = 0, \dots, d$ 
             $w_{ij} \leftarrow w_{ij} + \eta(r_i^t - y_i^t)x_j^t$ 
    Until 收敛
  
```

图 11-3 对于具有 $K > 2$ 个类的情况, 实现随机梯度下降的感知器训练算法。这是图 10-8 中给出的算法的在线版本

式(11.7)和式(11.9)都具有如下形式

$$\text{更新} = \text{学习因子} \times (\text{期望输出} - \text{实际输出}) \times \text{输入} \quad (11.10)$$

让我们更深入地考察上式。首先,如果实际输出等于期望输出,则不需要更新。当进行更新时,更新随期望输出与实际输出的差增加而增加。我们还看到,如果实际输出小于期望输出,则当输入为正时更新为正,输入为负时更新为负。这具有增加实际输出和降低差别的效果。如果实际输出大于期望输出,则当输入为正时更新为负,输入为负时更新为正;这就降低了实际输出,使得它更接近于期望输出。

在做更新时,更新量还依赖于输入。如果输入接近于 0,则它对实际输出的影响很小,因此其权重用一个较小的量更新。输入越大,其权重的更新也越大。

最后,更新量依赖于学习因子 η 。如果它太大,则更新过分依赖当前实例;就像系统只有短期记忆。如果该因子太小,则可能需要很多次更新才收敛。在 11.8.1 节,我们将讨论加快收敛的方法。

11.4 学习布尔函数

在布尔函数中,输入是二元的,并且如果对应的函数值为真则输出为 1,否则为 0。这样,它可以看作两类分类问题。作为一个例子,考虑学习 AND 两个输入,输入表和期望输出显示在表 11-1 中。实现 AND 的感知器和它的二维几何表示的一个例子显示在图 11-4 中。判别式是

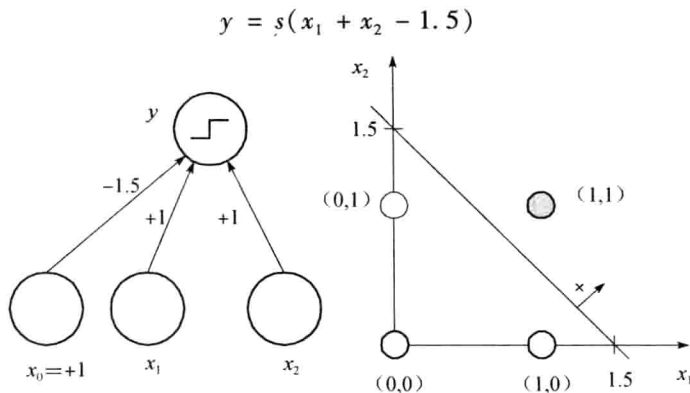


图 11-4 实现 AND 的感知器和它的几何表示

也就是说 $\mathbf{x} = [1, x_1, x_2]^T$, $\mathbf{w} = [-1.5, 1, 1]^T$ 。注意, $y = s(x_1 + x_2 - 1.5)$ 满足表 11-1 中 AND 函数定义给定的四个约束条件。例如,对于 $x_1 = 1, x_2 = 0$, $y = s(-0.5) = 0$ 。类似地,可以证明 $y = s(x_1 + x_2 - 0.5)$ 实现 OR。

表 11-1 AND 函数的输入和输出

x_1	x_2	r
0	0	0
0	1	0
1	0	0
1	1	1

尽管像 AND 和 OR 这样的布尔函数是线性可分的，并且是可以用感知器求解的，但是像 XOR 这样的函数不是。XOR 的输入和期望输出在表 11-2 中。正如我们可以从图 11-5 看到的，该问题不是线性可分的。可以证明这一点。注意不存在 w_0 , w_1 和 w_2 的值满足下列不等式：

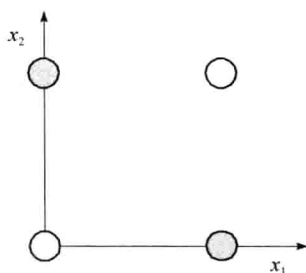
$$\begin{aligned} w_0 &\leq 0 \\ w_2 + w_0 &> 0 \\ w_1 + w_0 &> 0 \\ w_1 + w_2 + w_0 &\leq 0 \end{aligned}$$


图 11-5 XOR 问题不是线性可分的。我们不能划一条直线使得空心圆在一侧，实心圆在另一侧

表 11-2 XOR 函数的输入和输出

x_1	x_2	r
0	0	0
0	1	1
1	0	1
1	1	0

我们对这一结果并不奇怪，因为(二维)直线的 VC 维为三。具有二元输入，存在四种情况，因此我们知道存在具有两个输入的问题，它们不能用直线求解；XOR 就是其中之一。

11.5 多层感知器

具有单层权重的感知器只能逼近输入的线性函数，不能解决像 XOR 这样的问题，这些问题的判别式是非线性的。类似地，这种感知器也不能用于非线性回归。对于输入和输出层之间存在中间层或隐藏层(hidden layer)的前馈网络，就不存在这种局限性。如果用于分类，这种多层感知器(multilayer perceptrons, MLP)可以实现非线性判别式，而如果用于回归，可以逼近输入的非线性函数。

输入 \mathbf{x} 提供给输入层(包括偏倚)，“活性”向前传播，并计算隐藏单元的值 z_h (参见图 11-6)。每个隐藏单元自身都是一个感知器，并将非线性的 S 形函数作用于它的加权和：

244
245

$$z_h = \text{sigmoid}(\mathbf{w}_h^T \mathbf{x}) = \frac{1}{1 + \exp \left[- \left(\sum_{j=1}^d w_{hj} x_j + w_{h0} \right) \right]}, \quad h = 1, \dots, H \quad (11.11)$$

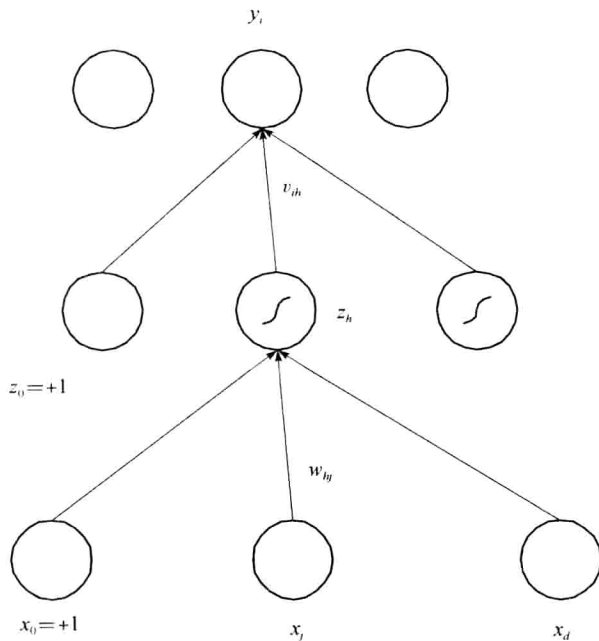


图 11-6 多层感知器结构。 $x_j (j=0, \dots, d)$ 是输入, $z_h (h=1, \dots, H)$ 是隐藏单元, 其中 H 是隐藏空间的维度。 z_0 是隐藏层的偏倚。 $y_i (i=1, \dots, K)$ 是输出单元。 w_{hj} 是第一层的权重, 而 v_{ih} 是第二层的权重

输出 y_i 是在第二层的感知器, 取隐藏单元作为它们的输入

$$y_i = \mathbf{v}_i^T \mathbf{z} = \sum_{h=1}^H v_{ih} z_h + v_{i0} \quad (11.12)$$

其中隐藏层还有一个偏倚单元, 记作 z_0 , 而 v_{i0} 是偏倚权重。输入层单元 x_j 不计, 因为计算不在那里进行, 并且当有一个隐藏层时, 网络是两层网络。

与通常一样, 在回归问题中, 在计算 y 的输出层不存在非线性。在两类判别式任务中, 有一个 S 形输出单元, 并且在存在 $K > 2$ 个类时, 有 K 个以软最大作为输出非线性的输出。

如果隐藏层单元的输出是线性的, 则隐藏层就没有用: 线性组合的线性组合还是一种线性组合。S 形函数是取阈值的连续、可微版本。我们需要可微性, 因为我们将看到学习方程是基于梯度的。另一种可以使用的 S 形非线性基函数是双曲正切函数 \tanh , 它值域是 -1 到 $+1$, 而不是 0 到 $+1$ 。在实践中, 使用 sigmoid 与使用 \tanh 并无区别。还有一种可能是使用高斯函数, 它使用欧氏距离而不是用点积表示相似性; 我们将在第 12 章讨论这种径向基函数网络。

输出是隐藏层单元计算的非线性基函数值的线性组合。可以说隐藏单元做了一个从 d -维输入空间到隐藏单元生成的 H -维空间的非线性变换, 并且在这个空间中, 输出层实现了一个线性函数。

我们不限于只有一个隐藏层, 而是可以将更多的、具有自己的输入权重的隐藏层放置在具有 S 形隐藏单元的第二个隐藏层之后, 从而计算第二个隐藏单元层的非线性函数, 实现输

人的更复杂的函数。实践中，人们很少构建超过一个隐藏层的网络，因为分析多个隐藏层的网络相当复杂。但是，有时隐藏层包含的隐藏单元太多时，使用多个隐藏层可能是明智的，宁可要“长而窄”的网络，也不要“短而胖”的网络。

11.6 作为普适近似的 MLP

我们可以将任意布尔函数表示成合取的析取，这样一个布尔表达式可以用一个具有一个隐藏层的多层感知器实现。每个合取用一个隐藏单元实现，而析取用输出单元实现。例如，

$$x_1 \text{ XOR } x_2 = (x_1 \text{ AND } \sim x_2) \text{ OR } (\sim x_1 \text{ AND } x_2)$$

前面我们已经看到如何使用感知器实现 AND 和 OR。因此，两个感知器可以平行地实现两个 AND，而另一个感知器可以将它们 OR 在一起(参见图 11-7)。我们看到，第一个隐藏层将输入从 (x_1, x_2) 映射到由第一层感知器定义的 (z_1, z_2) 空间。注意，输入 $(0, 0)$ 和 $(1, 1)$ 都被映射到 (z_1, z_2) 空间的 $(0, 0)$ ，使得在第二个空间是线性可分的。

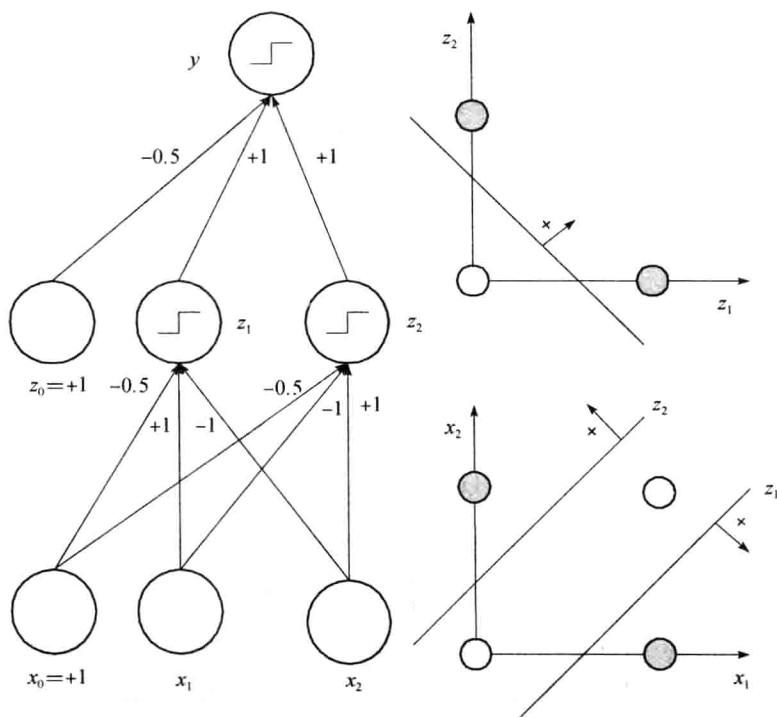


图 11-7 求解 XOR 问题的多层感知器。隐藏单元和输出单元具有阈值在 0 上的阈值激活函数

这样，在二元情况下，对于输出为 1 的每个输入组合，我们定义一个隐藏单元，它检查输入的这个特定合取。然后，输出单元实现析取。注意，这只是一个存在性证明，而这种网络可能不现实，因为当存在 d 个输入时，可能需要多达 2^d 个隐藏单元。这种结构实现了表查找而不是一般化。

我们可以将这些扩展到输入是连续值的情况，并且类似地证明具有连续输入和输出的任何函数都可以用多层感知器近似。使用两个隐藏层，普适近似(universal approximation)的证明很容易：对于每种输入或区域，使用第一个隐藏层上的隐藏单元，该区域可以被所有边上的超平

面所界定。第二个隐藏层的单元 AND 它们，围住该区域。然后，将隐藏单元到输出单元的连接权重设置为期望的函数值。这给出函数的分段常量近似 (piecewise constant approximation)；这对应忽略泰勒展开式中除常数项之外的所有项。增加隐藏单元的数量，并在输入空间中取更细的栅格，可以提高逼近期望值的精度。注意，没有给定期望的隐藏单元个数的形式上界。这种性质只是确保存在一个解；除此之外对我们并无其他帮助。业已证明，具有一个隐藏层的 MLP (具有任意个数的隐藏单元) 可以学习输入的任意非线性函数 (Hornik、Stinchcombe 和 White 1989)。

11.7 后向传播算法

训练多层感知器与训练感知器一样；唯一的区别是现在的输出是输入的非线性函数，这要感谢隐藏单元中的非线性偏倚函数。考虑把隐藏单元作为输入，第二层是感知器，我们已经知道在给定输入 z_h 的情况下，如何更新参数 v_{ij} 。对于第一层权重 w_{hj} ，我们使用链规则计算梯度：

$$\frac{\partial E}{\partial w_{hj}} = \frac{\partial E}{\partial y_i} \frac{\partial y_i}{\partial z_h} \frac{\partial z_h}{\partial w_{hj}}$$

误差就像从输出 y 传回到输出一样，因此新创了术语后向传播 (backpropagation) (Rumelhart、Hinton 和 Williams 1986a)。

11.7.1 非线性回归

让我们首先考虑用下式计算的 (具有单个输出的) 非线性回归：

$$y^t = \sum_{h=1}^H v_h z_h^t + v_0 \quad (11.13)$$

其中 z_h 用式 (11.11) 计算。在整个回归样本上的误差函数是

$$E(\mathbf{W}, \mathbf{v} | \mathcal{X}) = \frac{1}{2} \sum_t (r^t - y^t)^2 \quad (11.14)$$

第二层是以隐藏单元为输入的感知器，并且我们使用最小二乘规则来更新第二层的权重：

$$\Delta v_h = \eta \sum_t (r^t - y^t) z_h^t \quad (11.15)$$

第一层也是感知器，以隐藏单元作为输出单元，但在更新第一层权重时，我们不能直接使用最小二乘规则，因为对于这些隐藏单元，我们没有指定的期望输出。这正是链规则起作用的地方。我们有

$$\begin{aligned} \Delta w_{hj} &= -\eta \frac{\partial E}{\partial w_{hj}} \\ &= -\eta \sum_t \frac{\partial E^t}{\partial y^t} \frac{\partial y^t}{\partial z_h^t} \frac{\partial z_h^t}{\partial w_{hj}} \\ &= -\eta \sum_t \underbrace{-(r^t - y^t)}_{\partial E^t / \partial y^t} \underbrace{v_h}_{\partial y^t / \partial z_h^t} \underbrace{z_h^t (1 - z_h^t) x_j^t}_{\partial z_h^t / \partial w_{hj}} \\ &= \eta \sum_t (r^t - y^t) v_h z_h^t (1 - z_h^t) x_j^t \end{aligned} \quad (11.16)$$

前两项的乘积 $(r' - y')v_h$ 就像隐藏单元 h 的误差项。误差向后传播到隐藏单元。 $(r' - y')$ 是输出误差,按隐藏单元的“责任”加权,由其权重 v_h 给出。在第三项中, $z_h(1 - z_h)$ 是S形函数的导数, x'_j 是加权和关于权重 w_{hj} 的导数。注意,第一层权重的改变 Δw_{hj} 使用了第二层的权重 v_h 。因此,我们应当计算两层的改变,并更新第一层的权重,然后使用第二层权重的旧值更新第二层的权重。

最初,权重 w_{hj} 和 v_h 从小随机值(例如,区间 $[-0.01, 0.01]$ 中的值)开始,使得S形函数不饱和。规范化输入使得它们都具有均值0和单位方差并且具有相同尺度也是一种好的想法,因为我们使用了单个 η 参数。

使用这里给定的学习方程,对于每个模式,我们计算每个参数改变的方向和改变量。在批学习(batch learning),我们累积所有模式上的改变,并且在完全扫描了整个训练集之后做一次改变,如前面的更新方程所示。训练集中所有模式的一次完整扫描称作一个周期(epoch)。也可以在线学习,每个模式后更新权重,实现随机梯度下降。在这种情况下,应当选择较小的学习因子 η ,并且应当以随机次序扫描模式。因为数据集中可能有类似的模式,在线学习收敛较快,并且随机性具有增加噪声的效果,并有助于避免陷入局部极小。

为回归训练多层感知器的一个例子显示在图11-8中。随着训练继续,MLP拟合逐渐接近底层函数,并且误差降低(参见图11-9)。图11-10显示如何用隐藏单元输出的和形成MLP拟合。

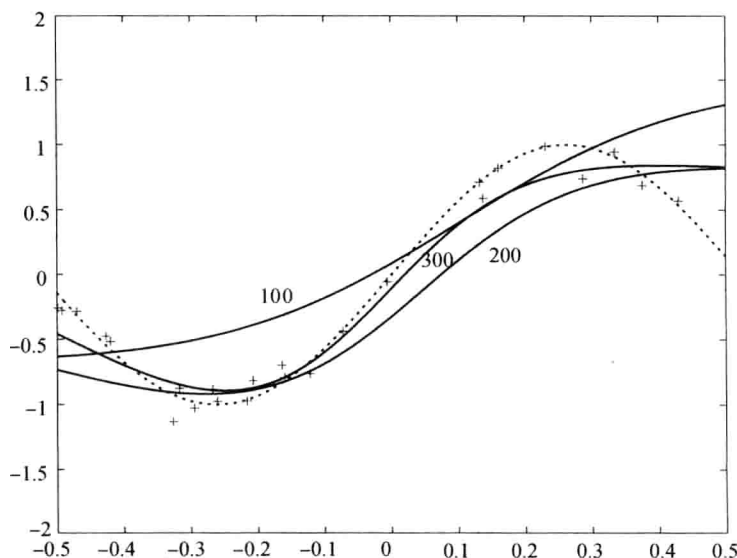


图 11-8 样本训练数据显示为“+”,其中 $x^i \sim U(-0.5, 0.5)$,而 $y^i = f(x^i) + \mathcal{N}(0, 0.1)$ 。 $f(x) = \sin(6x)$ 用虚线显示。图中绘制了100、200和300个周期后,具有两个隐藏单元的MLP的拟合演变

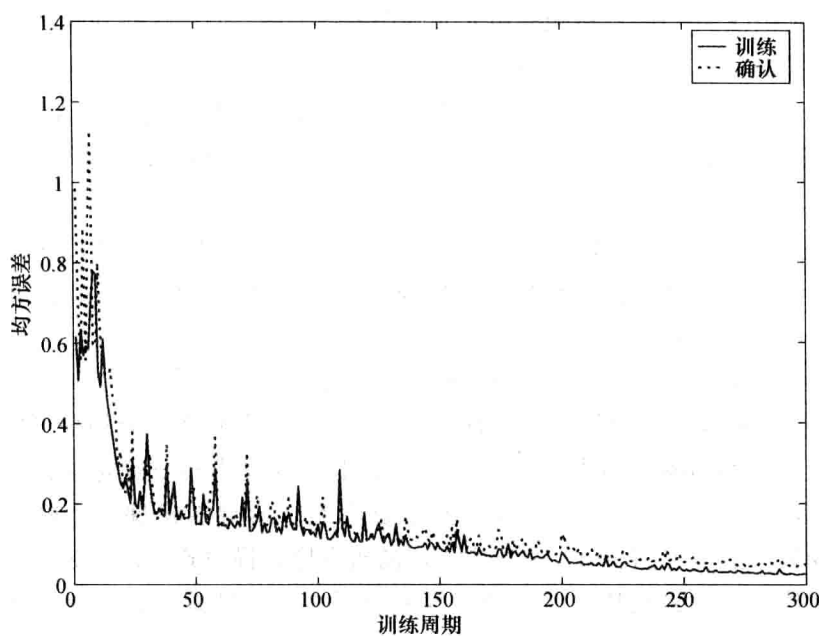


图 11-9 作为训练周期的函数，训练和验证集上的均方误差

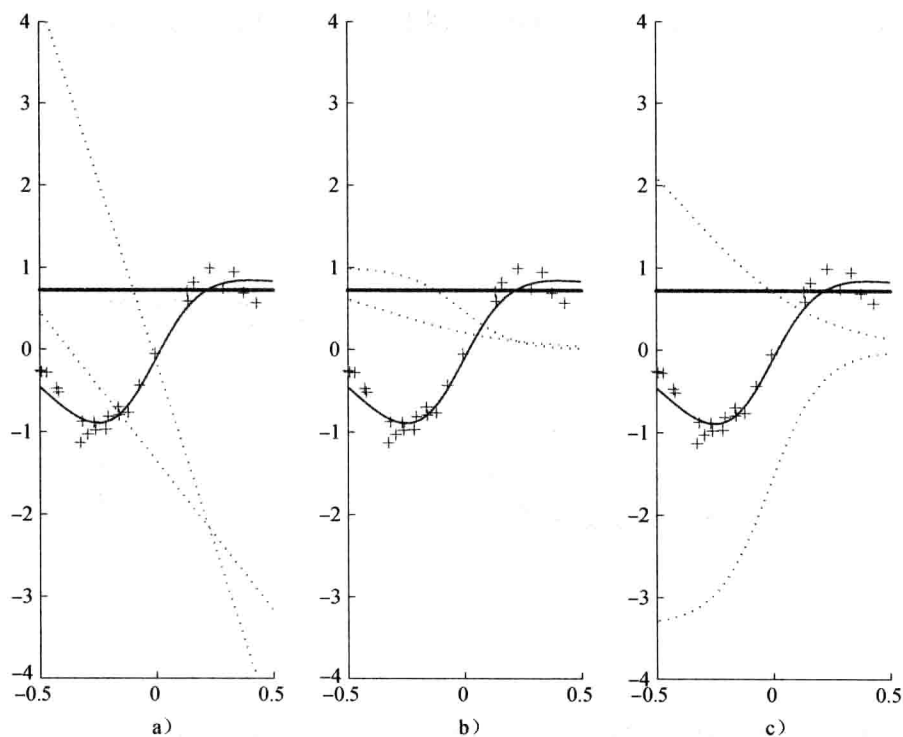


图 11-10 a) 第一层上隐藏单元权重的超平面, b) 隐藏单元输出, c) 隐藏单元输出乘以第二层的权重。纤细线显示的两个 S 形隐藏单元, 一个乘以负权重, 相加时实现隆起。使用更多的隐藏单元可以得到更好的近似 (参见图 11-12)

还可以有多个输出单元。在这种情况下,同时学习多个回归问题。我们有

$$y_i^t = \sum_{h=1}^H v_{ih} z_h^t + v_{i0} \quad (11.17)$$

而误差是

$$E(\mathbf{W}, \mathbf{V} | \mathcal{X}) = \frac{1}{2} \sum_t \sum_i (r_i^t - y_i^t)^2 \quad (11.18)$$

批更新规则为

$$\Delta v_{ih} = \eta \sum_t (r_i^t - y_i^t) z_h^t \quad (11.19)$$

$$\Delta w_{hj} = \eta \sum_t \left[\sum_i (r_i^t - y_i^t) v_{ih} \right] z_h^t (1 - z_h^t) x_j^t \quad (11.20)$$

$\sum_i (r_i^t - y_i^t) v_{ih}$ 是从所有输出单元累积的误差。算法的伪代码在图 11-11 中。注意,在这种情况下,所有输出单元共享相同的隐藏单元,因此使用相同的隐藏表示。一种可供选择的方法是对每个回归问题训练一个多层感知器,每个都有自己的隐藏单元。

将所有的 v_{ih} 和 w_{hj} 初始化为 $\text{rand}(-0.01, 0.01)$

Repeat

For 随机次序下所有的 $(\mathbf{x}^t, \mathbf{r}^t) \in \mathcal{X}$

For $h = 1, \dots, H$

$z_h \leftarrow \text{sigmoid}(\mathbf{w}_h^T \mathbf{x}^t)$

For $i = 1, \dots, K$

$y_i = \mathbf{v}_i^T \mathbf{z}$

For $i = 1, \dots, K$

$\Delta \mathbf{v}_i = \eta (r_i^t - y_i^t) \mathbf{z}$

For $h = 1, \dots, H$

$\Delta \mathbf{w}_h = \eta \left(\sum_i (r_i^t - y_i^t) v_{ih} \right) \mathbf{z}_h (1 - z_h) \mathbf{x}^t$

For $i = 1, \dots, K$

$\mathbf{v}_i \leftarrow \mathbf{v}_i + \Delta \mathbf{v}_i$

For $h = 1, \dots, H$

$\mathbf{w}_h \leftarrow \mathbf{w}_h + \Delta \mathbf{w}_h$

Until 收敛

图 11-11 为具有 K 个输出的回归训练多层感知器的后向传播算法。容易调整代码用于两类问题(设置单个 S 形输出)和 $K > 2$ 类问题(使用软最大输出)

11.7.2 两类判别式

在只有两个类时,一个输出单元就足够了:

$$y^t = \text{sigmoid} \left(\sum_{h=1}^H v_h z_h^t + v_0 \right) \quad (11.21)$$

它近似 $P(C_1 | \mathbf{x}^t)$ 和 $\hat{P}(C_2 | \mathbf{x}^t) \equiv 1 - y^t$ 。从 10.7 节中,我们知道在此情况下,误差函数是

$$E(\mathbf{W}, \mathbf{v} | \mathcal{X}) = - \sum_t r^t \log y^t + (1 - r^t) \log(1 - y^t) \quad (11.22)$$

实现梯度下降的更新方程是

$$\Delta v_h = \eta \sum_t (r^t - y^t) z_h^t \quad (11.23)$$

$$\Delta w_{hj} = \eta \sum_i (r^i - y^i) v_h z_h^i (1 - z_h^i) x_j^i \quad (11.24)$$

与简单感知器一样，回归和分类的更新方程是相同的(这不意味它们的值相同)。

11.7.3 多类判别式

在 $K > 2$ 类分类问题中，有 K 个输出

$$o_i^t = \sum_{h=1}^H v_{ih} z_h^t + v_{i0} \quad (11.25)$$

并且我们使用软最大指示类之间的依赖性，即它们是互斥的和穷举的：

$$y_i^t = \frac{\exp o_i^t}{\sum_k \exp o_k^t} \quad (11.26)$$

其中 y_i 近似 $P(C_i | \mathbf{x}^t)$ 。误差函数是

$$E(\mathbf{W}, \mathbf{V} | \mathcal{X}) = - \sum_i \sum_j r_i^t \log y_i^t \quad (11.27)$$

并且我们使用梯度下降得到更新方程：

$$\Delta v_{ih} = \eta \sum_i (r_i^t - y_i^t) z_h^t \quad (11.28)$$

$$\Delta w_{hj} = \eta \sum_i \left[\sum_i (r_i^t - y_i^t) v_{ih} \right] z_h^t (1 - z_h^t) x_j^t \quad (11.29)$$

Richard 和 Lippmann(1991) 证明，给定一个足够复杂的网络和足够的训练数据，适当训练的多层感知器可估计后验概率。

11.7.4 多个隐藏层

正如我们在前面看到的，可以有多个隐藏层，每个具有自己的权重，并将 S 形函数作用于它的加权和。对于回归，假设我们有一个多层感知器，具有两个隐藏层，我们有

$$z_{1h} = \text{sigmoid}(\mathbf{w}_{1h}^T \mathbf{x}) = \text{sigmoid} \left(\sum_{j=1}^d w_{1hj} x_j + w_{1h0} \right), h = 1, \dots, H_1$$

$$z_{2l} = \text{sigmoid}(\mathbf{w}_{2l}^T \mathbf{z}_1) = \text{sigmoid} \left(\sum_{h=0}^{H_1} w_{2lh} z_{1h} + w_{2l0} \right), l = 1, \dots, H_2$$

$$y = \mathbf{v}^T \mathbf{z}_2 = \sum_{l=1}^{H_2} v_l z_{2l} + v_0$$

其中 \mathbf{w}_{1h} 和 \mathbf{w}_{2l} 分别是第一和第二层的权重， z_{1h} 和 z_{2h} 分别是第一和第二个隐藏层的单元，而 \mathbf{v} 是第三层的权重。训练这种网络是类似的，唯一的区别在于，为了训练第一层的权重，我们需要后向传播更多层(习题 5)。

11.8 训练过程

11.8.1 改善收敛性

梯度下降具有多种优点。它简单。它是局部的，即权重的改变只使用前后突触单元和误

差(适合后向传播)的值。当使用在线训练时,它不需要存储训练集,并且可以自适应学习任务的变化。由于这些原因,它可以(并且已经)用硬件实现。但是,就自身而言,梯度下降收敛很慢。当学习时间很重要时,可以使用更复杂的优化方法(Battiti 1992)。Bishop (1995)详细讨论了多层感知器训练的共轭梯度的应用和二阶方法。然而,有两种频繁使用的简单技术,可以显著地改善梯度下降的性能,使得基于梯度的方法在实际应用中是可行的。

1. 动量

令 w_i 为多层感知器任意层中的任意权重,包括偏倚。在每次参数更新,相继的 Δw_i^t 可能很不相同以至于可能出现摆动,减缓收敛。 t 为时间指数,是批学习的周期数和在线学习的迭代次数。基本思想是在当前的改变中考虑上一次的更新,取移动平均,就好像因上次更新而存在动量(momentum):

$$\Delta w_i^t = -\eta \frac{\partial E^t}{\partial w_i} + \alpha \Delta w_i^{t-1} \quad (11.30)$$

通常, α 在 0.5 和 1.0 之间取值。当使用在线学习时,这种方法特别有用。我们将得到平均和光滑收敛轨迹的效果。缺点是需要将过去的 Δw_i^{t-1} 存放在附加的存储中。

2. 自适应学习率

在梯度下降中,学习因子 η 决定参数的改变量。它通常在 0.0 到 1.0 之间取值,大部分情况下小于或等于 0.2。为了更快收敛,可以让它自适应。学习进行时它保持较大,学习减慢时它也减小:

$$\Delta \eta = \begin{cases} +a & \text{如果 } E^{t+\tau} < E^t \\ -b\eta & \text{否则} \end{cases} \quad (11.31)$$

这样,如果训练集上的误差减小,则 η 增加一个常量;如果误差增大,则 η 减小。由于 E 可能从一个周期到另一个周期震荡,因此最好用过去几个周期的平均值作为 E^t 。

11.8.2 过分训练

具有 d 个输入、 H 个隐藏单元、 K 个输出的多层感知器的第一层有 $H(d+1)$ 个权重,第二层有 $K(H+1)$ 个权重。MLP 的时间和空间复杂度都是 $\mathcal{O}(H \cdot (K+d))$ 。用 e 表示训练周期数,则训练时间复杂度为 $\mathcal{O}(e \cdot H \cdot (K+d))$ 。

在一个应用中, d 和 K 是预先确定的, H 是参数,我们用它来调整模型的复杂性。从前面的章节中我们知道,过于复杂的模型记住了训练集中的噪声,不能泛化到验证集。例如,先前我们在多项式回归中已经看到这种现象,在那里我们看到噪声或小样本的出现增加了多项式的阶,导致更糟糕的泛化。类似地,在 MLP 中,当隐藏单元数很大时,泛化精度恶化(参见图 11-12),并且像任何统计学估计一样,对于 MLP,也存在偏倚/方差的两难选择(Geman、Bienenstock 和 Doursat 1992)。

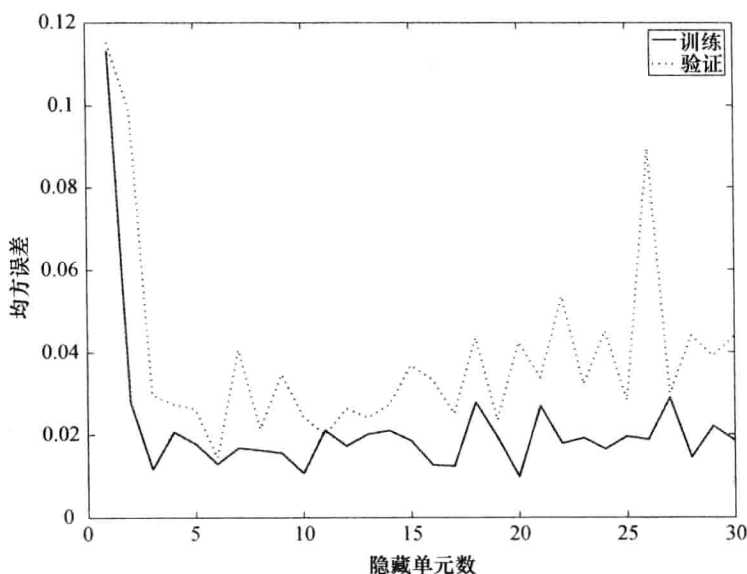


图 11-12 随着复杂度增加，训练误差固定但验证误差开始增加，网络开始过分拟合

当训练时间过长时，类似的事情也会发生：随着训练周期的增加，训练集上的误差降低，但是当超过某一点时，验证集上的误差开始增加(参见图 11-13)。回忆一下，最初所有的权重都接近于 0，因此影响很小。随着训练继续进行，大部分重要的权重开始离开 0 并发挥作用。但是，如果训练一直继续，得到训练集上越来越小的误差，几乎所有的权重都被更新，远离 0 成为有效的参数。这样，随着训练继续进行，就像将新的参数添加到系统中一样，增加了系统的复杂度，导致糟糕的泛化。学习应当在不是太晚时停止，以减轻过分训练 (overtraining) 问题。停止训练的最佳点和最佳隐藏单元数通过交叉验证确定，这涉及在训练期间未曾见过的验证集上测试网络的性能。

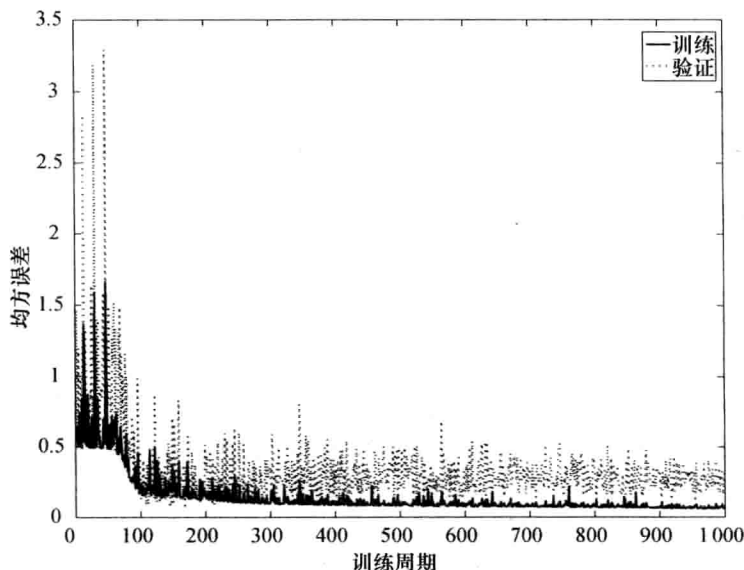


图 11-13 随着训练继续进行，验证误差开始增加，网络开始过分拟合

由于非线性性, 误差函数可能具有多个极小, 而梯度下降收敛于最近的极小。为了能够获得期望的误差, 通常以不同的初始权重开始, 对相同的网络训练多次, 并且计算验证误差的平均值。

11.8.3 构造网络

在某些应用中, 我们可能相信输入具有局部结构。例如, 在视频中, 我们知道邻近的像素是相关的, 并且存在诸如边、角等局部特征。任何对象, 如手写体数字, 都可以定义为这些图元的组合。类似地, 在语音中, 存在时间局部性, 并且时间上相近的输入可能组成语音基元。组合这些基元, 可以定义较长的话语, 如语音音素。在这种情况下, 在设计 MLP 时, 并不是将隐藏单元连接到所有的输入单元, 因为并非所有的输入都是相关的。替换地, 我们定义隐藏单元, 它在输入空间上定义一个输入窗口, 并且仅与输入的一个小的局部子集相连接。这样做减少了连接数, 从而减少了自由参数的数目 (Le Cun 等 1989)。

258

我们可以在相继层重复这一做法, 直到输出层。每层连接下一层的少量局部单元, 并且通过组合下面输入空间的较大部分, 检测更复杂的特征, 直到输出单元 (参见图 11-14)。例如, 输入可能是像素。通过观察像素, 第一个隐藏层的单元可以学习检测各方向的边。然后, 通过组合一些边, 第二个隐藏层的单元可以学习检测边的组合 (例如弧、角、线段), 并且在较高层组合它们。这些单元可以寻找半圆、矩形, 或者在脸部识别应用中, 寻找眼、嘴等。这是层次锥体 (hierarchical cone) 的一个例子, 随着我们沿着网络向上直到我们得到类, 特征越来越复杂、抽象, 并且数量越来越少。

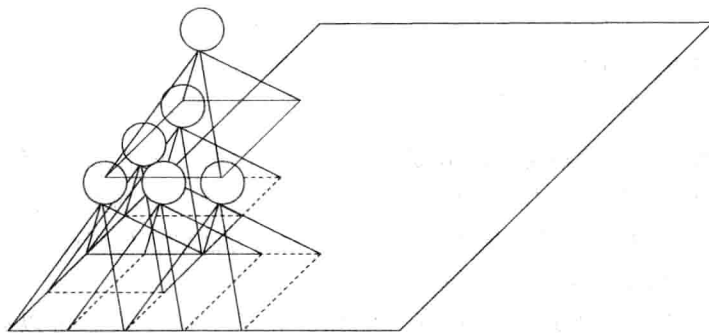


图 11-14 一个结构化的 MLP。每个单元都连接到其下单元的一个局部组群, 并检测一个特定的特征 (例如视频中的边、角等)。每个区域只显示了一个隐藏单元。通常, 有许多隐藏单元, 检测不同的局部特征

在这种情况下, 我们可以通过权重共享 (weight sharing) 进一步减少参数的数目。再次以视频识别为例, 我们可以看到: 在我们寻找类似有向边这样的特征时, 它们可能出现在输入空间的不同部分。因此, 替代定义独立的隐藏单元学习输入空间不同部分的不同特征, 我们可以有考察输入空间不同部分的相同隐藏单元的复制 (参见图 11-15)。在学习期间, 我们取不同的输入计算梯度, 然后对它们取平均值, 并做单个更新。这意味单个参数定义多个连接上的权重。此外, 由于一个权重上的更新基于多个输入的梯度, 因此训练集实际上就像有许

259
260

多个。

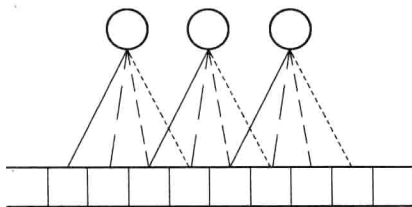


图 11-15 在权重共享中，不同的单元具有到不同输入的连接，但是共享相同的权重值（用线型表示）。只显示了一组单元；应当有多组单元，每个检测不同的特征

11.8.4 线索

局部结构的知识使得我们可以预先构造多层网络，并且使用权重共享，使得它具有较少的参数。具有全连接层的 MLP 不具有这种结构，并且更难训练。可能的话，与应用相关的任何类型的知识都应当构建到网络结构中。这些称作线索(hints)(Abu-Mostafa 1995)。它们是我们知道的目标函数的性质，独立于训练实例。

在图像识别，存在一些不变性线索：对象旋转、变换或缩放时，它的恒等性不变(参见图 11-16)。线索是辅助信息，可以用来指导学习过程，并且在训练集有限时特别有用。使用线索可以有不同的方法：

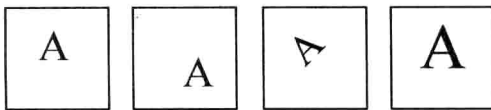


图 11-16 当对象变换、旋转或缩放时，它的恒等性不变。注意这并非总是为真，或者可能在某种程度为真：“b”和“q”互为旋转版本。这些是可以纳入到学习过程中使学习更容易的线索

1) 可以使用线索创建虚拟实例(virtual example)。例如，知道对象是缩放不变的，由给定的训练实例，我们可以用不同的尺寸产生多个复制，并以相同的类标号将它们添加到训练集中。这样做的优点是增大了训练集并且不需要对学习方法做任何修改。问题可能是对于学习方法，可能需要太多实例来学习不变性。

2) 不变性可以作为预处理阶段实现。例如，光学字符读入器可以有一个预处理步骤，将输入字符的图像关于尺寸和倾斜进行中心化和规范化。当可行时，这是最简单的解决方案。

3) 线索可以纳入到网络结构中。我们在 11.8.3 节看到的局部结构和权重共享就是一个例子，那里我们对小变换和旋转得到不变性。

4) 也可以通过修改误差函数纳入线索。假设我们知道从应用角度来说 \mathbf{x} 和 \mathbf{x}' 是相同的，其中 \mathbf{x}' 是 \mathbf{x} 的“虚拟实例”。也就是说，当 $f(\mathbf{x})$ 是我们逼近的函数时， $f(\mathbf{x}) = f(\mathbf{x}')$ 。让我们用 $g(\mathbf{x} | \theta)$ 表示我们的逼近函数，例如 MLP，其中 θ 是它的权重。然后，对于所有这样的 $(\mathbf{x}, \mathbf{x}')$ ，我们定义罚函数

$$E_h = [g(\mathbf{x} | \theta) - g(\mathbf{x}' | \theta)]^2$$

并把它作为一个额外项加到通常的误差函数中：

$$E' = E + \lambda_h \cdot E_h$$

这是一个罚项，处罚预测不服从线索的案例，而 λ_h 是这种罚的权重 (Abu-Mostafa 1995)。

另一个例子是近似线索：假设对于 x ，我们不知道准确的 $f(x)$ 值，但是我们知道它在区间 $[a_x, b_x]$ 中，则我们添加的罚项是：

$$E_h = \begin{cases} 0 & \text{如果 } g(x | \theta) \in [a_x, b_x] \\ (g(x) - a_x)^2 & \text{如果 } g(x | \theta) < a_x \\ (g(x) - b_x)^2 & \text{如果 } g(x | \theta) > b_x \end{cases}$$

这类似于支持向量机回归中使用的误差函数 (10.9.4 节)，它容忍小近似误差。

还有一个例子是正切支撑 (tangent prop) (Simard 等 1992)，其中变换与我们定义的线索相对。例如，旋转一个角度用一个函数建模。通常的误差函数被修改 (添加另一个项)，使得参数可以沿着这条变换线移动而不改变误差。

11.9 调整网络规模

前面，我们看到，当网络太大，具有太多自由参数时，泛化可能不好。为了找到最佳网络规模，最常用的方法是尝试不同的结构，在训练集上训练它们，并选择对验证集泛化最好的结构。另一种方法是将结构自适应 (structural adaptation) 合并到学习算法中。有两种方法可以做这件事：

1) 在破坏性 (destructive) 方法中，我们从一个大网络开始，逐步删除不必要的单元和/或连接。

2) 在建设性 (constructive) 方法中，我们从一个小网络开始，逐步增加改善性能的单元和/或连接。

一种破坏性方法是权衰减 (weight decay)，其基本思想是删除不必要的连接。理想情况下，为了能够确定一个单元或连接是否必要，我们需要使用它训练一次，不使用它训练一次，并检查独立的验证集上的误差之差。这种开销很大，因为这件事需要在单元/连接的所有组合上做。

假设如果一个连接的权重为 0，则没有使用它。我们给每个连接一个衰减到 0 的趋势，使得它除非为了降低误差而被明显地加强，否则就将消失。对于网络中的任意权重 w_i ，我们使用更新规则：

$$\Delta w_i = -\eta \frac{\partial E}{\partial w_i} - \lambda w_i \quad (11.32)$$

这等价于在具有一个附加的罚项的误差函数上做梯度下降，惩罚具有许多非零权重的网络：

$$E' = E + \frac{\lambda}{2} \sum_i w_i^2 \quad (11.33)$$

较简单的网络是较好的泛化器暗示我们通过增加一个罚项实现。注意，我们并不是说简单的网络总是比大网络好；我们是说如果我们有两个具有相同的训练误差的网络，则较简单的那个 (即具有较少权重的那个) 有较高的可能性更好地泛化到验证集上。

式 (11.32) 中第二项的效果像一个弹簧，将每个权重拉向 0。从一个接近于 0 的值开始，除非实际误差的梯度很大并导致更新，否则由于第二项，权重将逐渐衰减为 0。 λ 是参数，

决定训练集上误差和由于非零参数导致的复杂性的相对重要性, 因此决定衰减速度: 使用大的 λ , 无论训练误差多大, 权重将被拉向 0; 使用小的 λ , 对非零权重的罚不大。使用交叉验证对 λ 进行微调。

取代从大网络开始并剪去不必要的连接或单元, 我们也可以从小网络开始, 必要时添加单元和相关的连接(参见图 11-17)。在动态节点创建(dynamic node creation)(Ash 1989)中, 训练具有一个隐藏层和一个隐藏单元的网络, 收敛后如果误差仍然很高, 则添加一个单元。随机初始化新添加的单元的输入权重和输出权重并与先前存在的权重一起训练。先前存在的权重不再重新初始化, 并从先前的值开始。

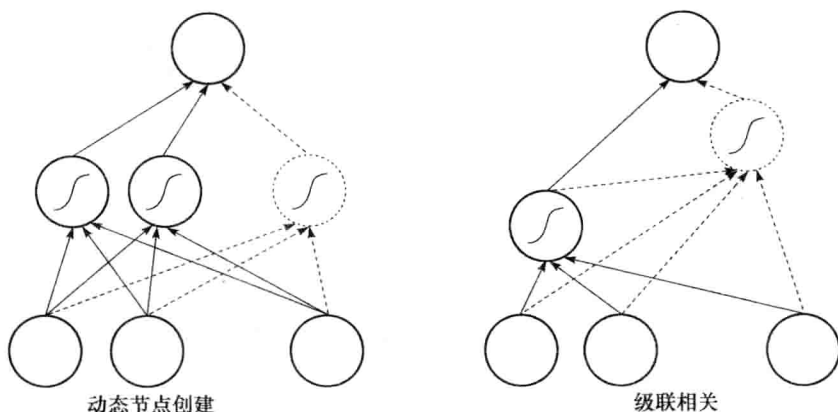


图 11-17 建设性方法的两个例子: 动态节点创建向一个已存在的隐藏层添加一个单元。级联相关添加每个单元作为新的隐藏层, 连接到前面的所有层。虚线表示新增加的单元/连接。为清晰起见, 忽略了偏倚单元/权重

在级联相关(cascade correlation)中(Fahlman 和 Lebiere 1990), 每个添加的单元是另一个隐藏层中的新的隐藏单元。每个隐藏层只有一个单元连接到它前面所有隐藏单元和输入。已存在的权重被冻结, 不再训练, 只训练新添加的单元的输入和输出权重。

动态节点创建在已经存在的隐藏层中创建一个新的隐藏单元, 而不增加新的隐藏层。级联关联总是创建具有单个单元的隐藏层。理想的建设性方法应当能够决定何时引进一个新的隐藏层, 何时向已有的隐藏层添加新单元。这是一个尚待解决的研究问题。

增量算法很有趣, 因为它在训练阶段不仅修改参数, 而且修改模型结构。我们可以考虑被多层感知器结构定义的空间和在该空间移动、对应于增加或删除单元以及层的操作(Aran 等 2009)。于是, 增量算法在这个状态空间搜索, (按照某种次序)尝试这些操作, 并根据某种优劣度量(例如复杂度和验证误差的某种组合)接受或拒绝。另一个例子是应用于多项式回归, 其中高阶项在训练阶段自动地添加/删除, 使得模型的复杂度与数据的复杂度相适应。随着计算费用逐渐降低, 这种自动的模型选择将成为学习过程的一部分自动地进行, 而不需要用户干预。

11.10 学习的贝叶斯观点

贝叶斯方法在训练神经网络时将参数(即连接权重 w_i)看作取自先验分布 $p(w_i)$ 的随机变

量, 并计算给定数据的后验概率

$$p(\mathbf{w} | \mathcal{X}) = \frac{p(\mathcal{X} | \mathbf{w})p(\mathbf{w})}{p(\mathcal{X})} \quad (11.34)$$

其中 \mathbf{w} 是网络的所有权重的向量。MAP 估计 $\hat{\mathbf{w}}$ 为后验的众数

$$\hat{\mathbf{w}}_{MAP} = \arg \max_{\mathbf{w}} \log p(\mathbf{w} | \mathcal{X}) \quad (11.35)$$

取式(11.34)的对数, 我们得到

$$\log p(\mathbf{w} | \mathcal{X}) = \log p(\mathcal{X} | \mathbf{w}) + \log p(\mathbf{w}) + C$$

右边的第一项是对数似然, 而第二项是先验概率的对数。如果权重是独立的, 并且先验概率取作高斯分布 $\mathcal{N}(0, 1/2\lambda)$

$$p(\mathbf{w}) = \prod_i p(w_i) \quad \text{其中 } p(w_i) = c \cdot \exp \left[-\frac{w_i^2}{2(1/2\lambda)} \right] \quad (11.36)$$

则 MAP 估计最小化增广的误差函数

$$E' = E + \lambda \|\mathbf{w}\|^2 \quad (11.37)$$

其中 E 为通常的分类或回归误差(负的对数似然)。这个增广误差正是我们在权重衰减(参见式(11.33))中使用的误差函数。使用较大的 λ 意味较小的参数可变性, 对它们施加更大的力量, 使之接近于 0, 并且更多地考虑先验而不是数据; 如果 λ 较小, 则允许较大的参数可变性。这种删除不必要的参数的方法在统计学中称作岭回归(ridge regression)。

这是使用代价函数、结合对数据的拟合和模型复杂度正则化(regularization)的另一个例子

$$\text{代价} = \text{数据错拟合} + \lambda \cdot \text{复杂度} \quad (11.38)$$

MacKay(1992a, b)讨论了在训练多层感知器时使用贝叶斯估计。我们将在第 14 章更详细地讨论贝叶斯估计。 266

经验表明, 训练后, 多层感知器的大部分权重都围绕 0 正态分布, 证明使用权重衰减是正确的。但是, 也并非总是这种情况。Nowlan 和 Hinton(1992)提出了软权重共享(soft weight sharing), 其中权重取自混合高斯分布, 允许它们形成多个而不是一个簇。此外, 这些簇的中心可以在任何地方, 而不必在 0, 并且具有可以修改的方差。这将式(11.36)的先验概率改变成 $M \geq 2$ 个高斯混合

$$p(w_i) = \sum_{j=1}^M \alpha_j p_j(w_i) \quad (11.39)$$

其中 α_j 是优先权, $p_j(w_i) \sim \mathcal{N}(m_j, s_j^2)$ 是高斯分量。 M 由用户设置, 而 α_j 、 m_j 和 s_j 从数据中学习。在训练阶段使用这种先验并用它的对数增广误差函数, 权重收敛以降低误差, 并且还自动地分组以提高对数先验。

11.11 维度归约

在多层感知器中, 如果隐藏单元数小于输入数, 则第一层执行维度归约。这种归约形式和隐藏单元生成的新空间依赖于 MLP 的训练目的。如果 MLP 用来分类, 输出单元紧随隐藏层, 则定义了新空间并且学习该映射, 以降低分类误差(参见图 11-18)。

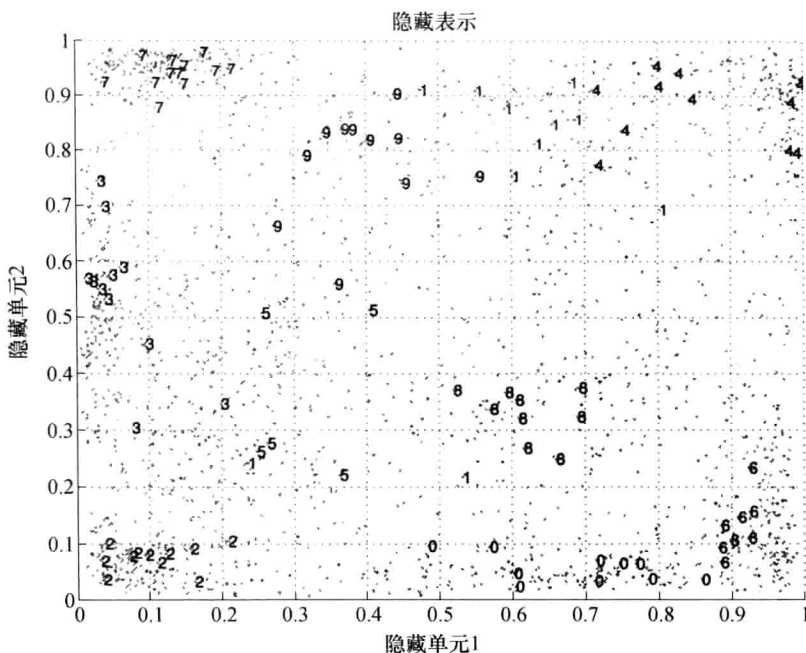


图 11-18 绘制在用于分类的训练后的 MLP 的两个隐藏单元的空间中的 Optdigits 数据。只显示了 100 个数据点的标号。该 MLP 具有 64 个输入，2 个隐藏单元和 10 个输出，具有 80% 的准确率。由于 S 形函数，隐藏单元的值在 0 和 1 之间，并且类在角落附近聚集。可以将该图与第 6 章的图比较。第 6 章的图在相同的数据集上使用其他维度归约方法绘制

通过分析权重，我们可以明白 MLP 是做什么的。我们知道当两个向量相等时点积最大。因此，我们可以认为每个隐藏单元定义了其输入权重的模板，并通过分析这些模板，我们可以从训练后的 MLP 中提取知识。如果输入是规范化的，则权重告诉我们它们的相对重要性。这样的分析并不容易，但是让我们洞察到 MLP 在做什么，并使得我们可以窥视黑箱。

一种有趣的结构是自动关联器 (autoassociator) (Cottrell、Munro 和 Zipser 1987)。这是一种 MLP 结构，其中输出与输入一样多，并且所要求的输出被定义等于输入 (参见图 11-19)。

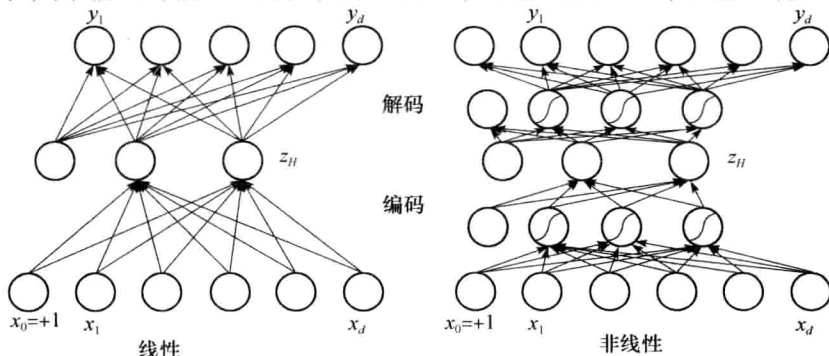


图 11-19 在自动关联器中，输出与输入一样多并且期望的输出是输入。当隐藏单元的数目小于输入的数目时，则 MLP 被训练以发现输入在隐藏层上的最佳编码，实现维度归约。左边，第一层充当编码器，而第二层充当解码器。右边，如果编码器和解码器是具有 S 形隐藏单元的多层感知器，则网络进行非线性维度归约

为了能够在输出层重新产生输入, MLP 被迫找出输入在隐藏层的最佳表示。当隐藏单元数小于输入数时, 这意味着维度归约。一旦训练完成, 从输入到隐藏层的第一层充当编码器, 而隐藏单元的值形成编码表示。从隐藏单元到输出单元的第二层充当解码器, 由原信号的编码表示重构原信号。

已经证明(Bourlard 和 Kamp 1988), 具有一个隐藏层的 MLP 实现主成分分析(6.3 节), 不同之处在于隐藏单元的权重不是使用特征值按重要性排序的特征向量, 但是它与 H 个主要特征向量生成相同的空间。如果编码器和解码器不是一层, 而是在隐藏单元具有 S 形非线性多层感知器, 则编码器实现非线性的维度归约。

另一种使用 MLP 进行维度归约的方法是通过多维定标(6.5 节)。Mao 和 Jain(1995)展示了如何使用 MLP 学习 Sammon 映射(Sammon mapping)。回忆式(6.29), Sammon 应力定义为

$$E(\theta | X) = \sum_{r,s} \left[\frac{\|g(\mathbf{x}^r | \theta) - g(\mathbf{x}^s | \theta)\| - \|\mathbf{x}^r - \mathbf{x}^s\|}{\|\mathbf{x}^r - \mathbf{x}^s\|} \right]^2 \quad (11.40)$$

一个具有 d 个输入, H 个隐藏单元和 $k < d$ 个输出单元的 MLP 用来实现 $g(\mathbf{x} | \theta)$, 将 d -维输入映射到一个 k -维向量, 其中 θ 对应于 MLP 的权重。给定一个数据集 $X = \{\mathbf{x}^i\}_i$, 我们可以使用梯度下降直接最小化 Sammon 应力来学习 MLP(即 $g(\mathbf{x} | \theta)$), 使得 k 维表示之间的距离与原空间中的距离尽可能接近。

11.12 学习时间

到目前为止, 我们一直关注输入一次全部提供的情况。在某些应用中, 输入是时间数据, 我们需要学习时间序列。换句话说, 输出也可能随时间变化。例子有

- 序列识别(sequence recognition)。这是把给定的序列指派到多个类中的一个。语音识别是一个例子, 其中输入信号序列是口语语音, 而输出是词的编码。即输入随时间变化, 但输出不随时间变化。
- 序列复制(sequence reproduction)。这里, 在看到给定序列的一部分之后, 系统将预测其余部分。时间序列产生是一个例子, 那里输入是给定的, 但输出是变化的。
- 时间关联(temporal association)。这是最一般的情况, 其中特定的输出序列作为特定的输入序列之后的输出。输入和输出序列可能不同。这里, 输入和输出都随时间变化。

11.12.1 时间延迟神经网络

识别时间序列的最简单的方法是把它转换成空间序列。然后可以利用前面讨论的任意方法进行分类。在时间延迟神经网络(time delay neural network)中(Waibel 等 1989), 前面的输入被延迟, 以便与最后的输入同步, 一起作为输入提交系统(参见图 11-20)。然后, 使用后向传播训练权重。为了提取局部于时间的特征, 我们可以使用结构化的连接层和权重共享, 以便得到时间的变换不变性。这种结构的主要限制是我们滑过的序列的时间窗口大小应当预先固定。

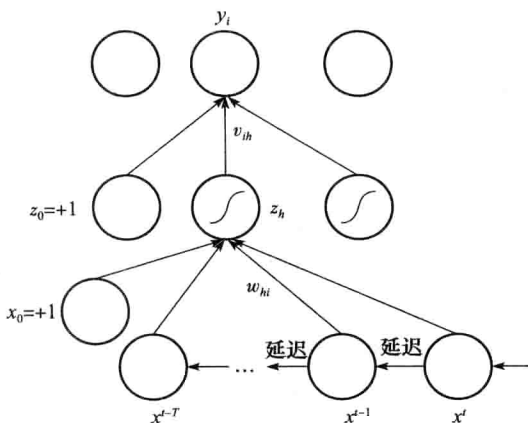


图 11-20 一个时间延迟网络。长度为 T 的时间窗口中的输入被延迟，直到我们可以将所有 T 个输入作为输入向量提供给 MLP

11.12.2 递归网络

在递归网络 (recurrent network) 中，除了前馈连接之外，单元具有自连接或到前面层的连接。这种递归性充当短期记忆，并使得网络记住过去发生的事。

在大部分情况下，我们使用部分递归网络，其中有限多个递归连接被添加到多层感知器中 (参见图 11-21)。这结合了多层感知器的非线性逼近能力和递归的时间表达能力的优点，并且这样的网络可以用来实现三种时间关联任务中的任何一种。还可以在递归的后向连接中

271

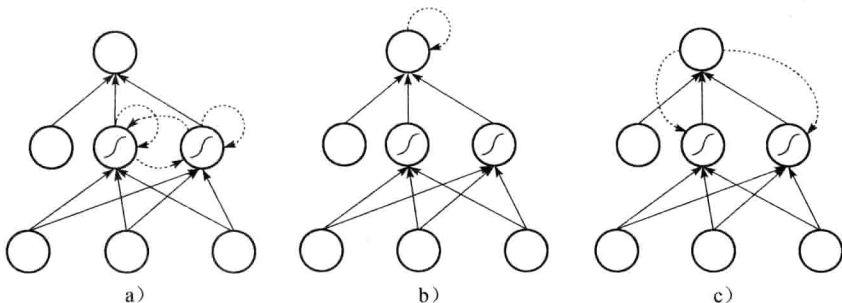


图 11-21 具有部分递归的 MLP 的例子。递归连接用虚线显示：a) 隐藏层中的自连接，b) 输出层中的自连接，c) 从输出层到隐藏层的连接。还可以有这些情况的组合

如果序列具有较小的最大长度，则可以使用按时间展开 (unfolding in time)，将任意的递归网络转换成等价的前馈网络 (参见图 11-22)。为不同时间的复制创建单独的单元和连接。结果网络可以用后向传播训练，附加的要求是每个连接的所有复制应当保持相等。与权重共享一样，方法是按时间对不同权重的改变求和，并用平均值更新权重。这称作通过时间后向传播 (backpropagation through time) (Rumelhart、Hinton 和 Williams 1986b)。这种方法的问题是如果序列的长度很长，则存储需求量很大。实时递归学习 (real time recurrent learning)

(Williams 和 Zipser 1989) 是一种训练递归网络而不展开的算法, 并且具有可以用于任意长度序列的优点。

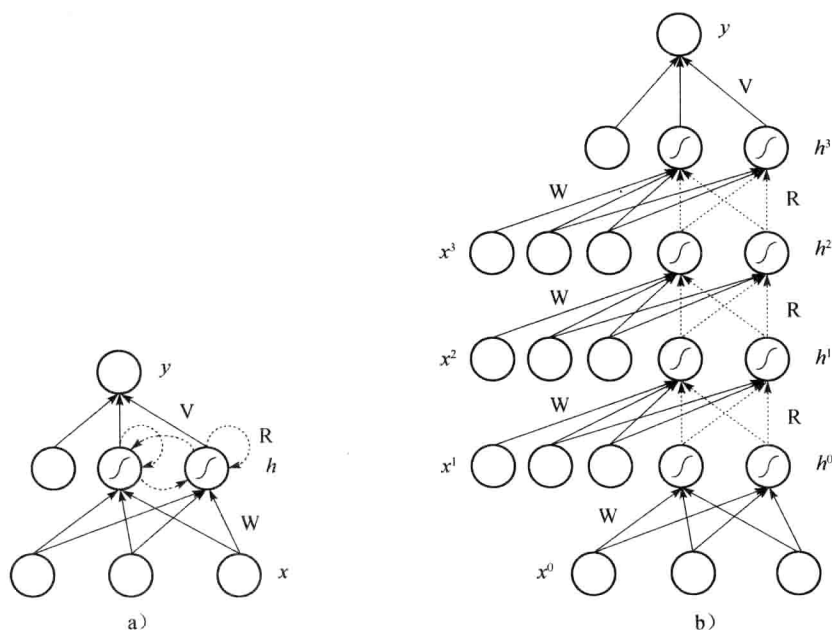


图 11-22 通过时间后向传播: a) 递归网络, b) 4 步中行为相同的、等价的展开网络

11.13 注释

人工神经网络的研究历史与数字计算机一样长。McCulloch 和 Pitts (1943) 提出了人工神经网络的第一个数学模型。Rosenblatt (1962) 提出了感知器模型和学习算法。Minsky 和 Papert (1969) 指出了单层感知器的局限性 (例如 XOR 问题), 并且由于那时还没有训练具有隐藏层的多层感知器的算法, 因此除了少数地方之外, 人工神经网络的工作几乎停止了。Hopfield (1982) 的文章带来了神经网络的复兴。随后出现了并行分布处理 (PDP) 研究小组编写的两卷并行处理书 (Rumelhart 和 McClelland 1986)。看起来, 后向传播几乎同时在多个地方被发明, 而单层感知器的局限性也不复存在。

从 20 世纪 80 年代中期开始, 出现了关于人工神经网络模型的大量研究, 来自各个学科: 物理学、统计学、心理学、认知科学、神经系统科学、语言学, 更不必说计算机科学、电子学和自适应控制了。或许, 人工神经网络研究的最重要贡献是这种沟通不同学科, 特别是沟通统计学与工程的协同。感谢这种协同, 使机器学习领域现在得以确立。

现在, 该领域更加成熟, 目标被更适当、更好地确定。对后向传播的批评之一是, 从生物学角度讲, 它几乎是不可能的! 尽管术语“神经网络”仍然被广泛使用, 但是通常把神经网络模型理解为 (例如多层感知器) 非参数估计方法, 并且分析它的最佳方法是使用统计学方法。

例如, 一种类似于多层感知器的统计学方法是投影追踪 (projection pursuit) (Friedman 和 Stuetzle 1981), 表示为

$$y = \sum_{h=1}^H \phi_h(\mathbf{w}_h^T \mathbf{x})$$

不同的是每个“隐藏单元”都具有自己的 $\phi_h(\cdot)$ ，尽管在 MLP 中，它们都是 S 型函数。在第 12 章，我们将看到另一种称作径向基函数的神经网络结构，它在隐藏单元使用高斯函数。

有各种各样的神经网络教科书：最早的是 Hertz、Krogh 和 Palmer 1991，仍然可以阅读。Bishop 1995 重点是模式识别并且详细讨论了可以用于训练的各种优化算法，以及贝叶斯方法，推广了权重衰减。Ripley 1996 从统计学的角度分析了神经网络。

人工神经网络，例如多层感知器，具有各种各样的成功应用。除了成功地用于自适应控制、语音识别和视频之外，有两点值得注意：Tesauro 的 TD-Gammon 程序 (Tesauro 1994) 使用增强学习 (第 16 章) 训练多层感知器，并在大师级玩西洋双陆棋。Pomerleau 的 ALVINN 是一个神经网络，通过观察驾驶员学习 5 分钟之后，它可以自动驾驶厢式货车，速度达每小时 20 英里。

11.14 习题

1. 给出一个可以计算其输入的 NOT 的感知器。
2. 给出一个可以计算其两个输入的 NAND 的感知器。
3. 给出一个可以计算其三个输入的奇偶性的感知器。
4. 当隐藏单元使用 tanh 函数而不是使用 S 形函数时，推导出更新方程。使用事实 $\tanh' = (1 - \tanh^2)$ 。
5. 为具有两个隐藏层的 MLP 推导更新方程。
6. 考虑一个具有一个隐藏层的 MLP 结构，其中还存在直接从输入到输出单元的权重。解释这种结构何时是有益的，如何训练它。
7. 奇偶性是循环移动不变的。例如，“0101”和“1010”具有相同的奇偶性。使用这个提示，提出一个学习奇偶函数的多层感知器。
8. 在级联相关中，冻结前面已经存在的权重有何优点？
9. 为实现最小化 Sammon 应力式 (11.40) 的 Sammon 映射的 MLP 推导更新方程。
10. 11.6 节我们讨论了一个具有两个隐藏层的 MLP 如何实现分段常数近似。证明：如果最后一层的权重不是常数而是输入的线性函数，则我们可以实现分段线性近似。
11. 为软权重共享推导更新方程。
12. 在自动关联网中，如何决定隐藏单元的个数？
13. MLP 结构的增量学习可以看作状态空间搜索。操作是什么？优度函数是什么？什么类型的搜索策略是合适的？以这种方式定义，使得动态结点创建和级联相关都是特例。
14. 对图 11.22 给出的 MLP，为展开网络推导更新方程。

11.15 参考文献

Abu-Mostafa, Y. 1995. “Hints.” *Neural Computation* 7: 639–671.

- Aran, O., O. T. Yıldız, and E. Alpaydın. 2009. "An Incremental Framework Based on Cross-Validation for Estimating the Architecture of a Multilayer Perceptron." *International Journal of Pattern Recognition and Artificial Intelligence* 23: 159-190.
- Ash, T. 1989. "Dynamic Node Creation in Backpropagation Networks." *Connection Science* 1: 365-375.
- Battiti, R. 1992. "First- and Second-Order Methods for Learning: Between Steepest Descent and Newton's Method." *Neural Computation* 4: 141-166.
- Bishop, C. M. 1995. *Neural Networks for Pattern Recognition*. Oxford: Oxford University Press.
- Bourlard, H., and Y. Kamp. 1988. "Auto-Association by Multilayer Perceptrons and Singular Value Decomposition." *Biological Cybernetics* 59: 291-294.
- Cottrell, G. W., P. Munro, and D. Zipser. 1987. "Learning Internal Representations from Gray-Scale Images: An Example of Extensional Programming." In *Ninth Annual Conference of the Cognitive Science Society*, 462-473. Hillsdale, NJ: Erlbaum.
- Durbin, R., and D. E. Rumelhart. 1989. "Product Units: A Computationally Powerful and Biologically Plausible Extension to Backpropagation Networks." *Neural Computation* 1: 133-142.
- Fahlman, S. E., and C. Lebiere. 1990. "The Cascade Correlation Architecture." In *Advances in Neural Information Processing Systems 2*, ed. D. S. Touretzky, 524-532. San Francisco: Morgan Kaufmann.
- Friedman, J. H., and W. Stuetzle. 1981. "Projection Pursuit Regression." *Journal of the American Statistical Association* 76: 817-823.
- Geman, S., E. Bienenstock, and R. Doursat. 1992. "Neural Networks and the Bias/Variance Dilemma." *Neural Computation* 4: 1-58.
- Hertz, J., A. Krogh, and R. G. Palmer. 1991. *Introduction to the Theory of Neural Computation*. Reading, MA: Addison Wesley.
- Hinton, G. E., and R. R. Salakhutdinov. 2006. "Reducing the dimensionality of data with neural networks." *Science* 313: 504-507.
- Hopfield, J. J. 1982. "Neural Networks and Physical Systems with Emergent Collective Computational Abilities." *Proceedings of the National Academy of Sciences USA* 79: 2554-2558.
- Hornik, K., M. Stinchcombe, and H. White. 1989. "Multilayer Feedforward Networks Are Universal Approximators." *Neural Networks* 2: 359-366.
- Le Cun, Y., B. Boser, J. S. Denker, D. Henderson, R. E. Howard, W. Hubbard, and L. D. Jackel. 1989. "Backpropagation Applied to Handwritten Zipcode Recognition." *Neural Computation* 1: 541-551.
- MacKay, D. J. C. 1992a. "Bayesian Interpolation." *Neural Computation* 4: 415-447.
- MacKay, D. J. C. 1992b. "A Practical Bayesian Framework for Backpropagation Networks" *Neural Computation* 4: 448-472.
- Mao, J., and A. K. Jain. 1995. "Artificial Neural Networks for Feature Extraction and Multivariate Data Projection." *IEEE Transactions on Neural Networks* 6: 296-317.

- Marr, D. 1982. *Vision*. New York: Freeman.
- McCulloch, W. S., and W. Pitts. 1943. "A Logical Calculus of the Ideas Immanent in Nervous Activity." *Bulletin of Mathematical Biophysics* 5: 115-133.
- Minsky, M. L., and S. A. Papert. 1969. *Perceptrons*. Cambridge, MA: MIT Press. (Expanded ed. 1990.)
- Nowlan, S. J., and G. E. Hinton. 1992. "Simplifying Neural Networks by Soft Weight Sharing." *Neural Computation* 4: 473-493.
- Pomerleau, D. A. 1991. "Efficient Training of Artificial Neural Networks for Autonomous Navigation." *Neural Computation* 3: 88-97.
- Posner, M. I., ed. 1989. *Foundations of Cognitive Science*. Cambridge, MA: MIT Press.
- Richard, M. D., and R. P. Lippmann. 1991. "Neural Network Classifiers Estimate Bayesian *a Posteriori* Probabilities." *Neural Computation* 3: 461-483.
- Ripley, B. D. 1996. *Pattern Recognition and Neural Networks*. Cambridge, UK: Cambridge University Press.
- Rosenblatt, F. 1962. *Principles of Neurodynamics: Perceptrons and the Theory of Brain Mechanisms*. New York: Spartan.
- Rumelhart, D. E., G. E. Hinton, and R. J. Williams. 1986a. "Learning Representations by Backpropagating Errors." *Nature* 323: 533-536.
- Rumelhart, D. E., G. E. Hinton, and R. J. Williams. 1986b. "Learning Internal Representations by Error Propagation." In *Parallel Distributed Processing*, ed. D. E. Rumelhart, J. L. McClelland, and the PDP Research Group, 318-362. Cambridge, MA: MIT Press.
- Rumelhart, D. E., J. L. McClelland, and the PDP Research Group, eds. 1986. *Parallel Distributed Processing*. Cambridge, MA: MIT Press.
- Simard, P., B. Victorri, Y. Le Cun, and J. Denker. 1992. "Tangent Prop: A Formalism for Specifying Selected Invariances in an Adaptive Network." In *Advances in Neural Information Processing Systems 4*, ed. J. E. Moody, S. J. Hanson, and R. P. Lippman, 895-903. San Francisco: Morgan Kaufmann.
- Tesauro, G. 1994. "TD-Gammon, A Self-Teaching Backgammon Program, Achieves Master-Level Play." *Neural Computation* 6: 215-219.
- Thagard, P. 2005. *Mind: Introduction to Cognitive Science*. 2nd ed. Cambridge, MA: MIT Press.
- Waibel, A., T. Hanazawa, G. Hinton, K. Shikano, and K. Lang. 1989. "Phoneme Recognition Using Time-Delay Neural Networks." *IEEE Transactions on Acoustics, Speech, and Signal Processing* 37: 328-339.
- Williams, R. J., and D. Zipser. 1989. "A Learning Algorithm for Continually Running Fully Recurrent Neural Networks." *Neural Computation* 1: 270-280.

局部模型

我们继续讨论多层神经网络，考察第一层包含局部接受单元的模型；这些局部接受单元响应输入空间局部区域中的实例。上面第二层对这些局部区域学习回归或分类函数。我们讨论找出重要局部区域，以及这些区域中的模型的学习方法。

12.1 引言

进行函数逼近的一种方法是将输入空间划分成局部小片，并且在每个局部小片中分别学习拟合。在第7章，我们讨论了聚类的统计学方法，使得我们能够对输入实例分组并对输入分布建模。竞争方法是用于在线聚类的神经网络方法。本章，我们讨论 k -均值的在线版本以及两种神经网络扩展：自适应共鸣理论 (ART) 和自组织映射 (SOM)。

然后，我们讨论一旦输入局部化，如何实现监督学习。如果局部小片上的拟合是常量，则该技术称作径向基函数 (RBF) 网络；如果拟合是输入的线性函数，则称作混合专家技术 (MoE)。我们讨论回归和分类，并与第11章讨论的 MLP 方法进行比较。

279

12.2 竞争学习

在第7章，我们使用半参数高斯混合密度，假定输入来自 k 个高斯源中的一个。本节，我们做相同的假设，数据中存在 k 个分组 (或簇)，但是我们的方法不是概率方法，因为我们不将参数模型强加在数据源上。另一个区别是我们提出的学习方法是在线的：在训练阶段我们并没有全部样本；我们逐个接收实例并更新模型参数。使用术语竞争学习 (competitive learning) 是因为这些分组，更确切地说，代表这些分组的单元为成为代表实例而相互竞争。这种方法也称胜者全取 (winner-take-all)；就像一个分组获胜并得到更新，而其他分组则完全不更新一样。

与第7章讨论的方法相反，这些方法本身可以用于在线聚类。在线方法具有通常的优点：(1) 不需要附加的存储保存整个训练集；(2) 每步更新简单、易于实现 (例如用硬件实现)；(3) 输入的分布可以随时间而改变，并且模型可以自动地适应这些改变。如果我们使用批处理算法，我们将需要收集新样本，并且从头开始在整个样本上运行批处理方法。

从12.3节开始，我们还将讨论这种方法如何后跟一种监督方法，来学习回归和分类问题。这将是两阶段系统，可以用两层网络实现，其中第一阶段 (层) 对输入密度建模并找到相应的局部模型，而第二阶段是产生最终输出的局部模型。

12.2.1 在线 k -均值

在式(7.3)中，我们定义重构误差为

$$E(\{\mathbf{m}_i\}_{i=1}^k | \mathcal{X}) = \frac{1}{2} \sum_t \sum_i b_i^t \|\mathbf{x}^t - \mathbf{m}_i\|^2 \quad (12.1)$$

其中

$$b_i^t = \begin{cases} 1 & \text{如果 } \|x^t - m_i\| = \min_l \|x^t - m_l\| \\ 0 & \text{否则} \end{cases} \quad (12.2)$$

$X = \{x^t\}_t$ 是样本, 而 $m_i (i=1, \dots, k)$ 是簇中心。如果 m_i 是 x^t 的欧氏距离最接近的中心, 则 b_i^t 为 1。好像是所有的 $m_l (l=1, \dots, k)$ 竞争, 而 m_i 赢得竞争, 因为它是最近的。

k -均值的批处理算法按下式更新中心

$$m_i = \frac{\sum_t b_i^t x^t}{\sum_t b_i^t} \quad (12.3)$$

一旦使用式(12.2)选取获胜者, 它将最小化式(12.1)。正如我们先前看到的, 计算 b_i^t 和更新 m_i 的两个步骤迭代, 直到收敛。

通过进行随机梯度下降、逐个考虑实例、并在每一步进行少许更新而不忘记先前的更新, 我们可以得到在线 k -均值 (online k -means)。对于单个实例, 重构误差为

$$E^t(\{m_i\}_{i=1}^k | x^t) = \frac{1}{2} \sum_i b_i^t \|x^t - m_i\|^2 = \frac{1}{2} \sum_i \sum_{j=1}^d b_i^t (x_j^t - m_{ij})^2 \quad (12.4)$$

其中 b_i^t 的定义同式(12.2)。对上式使用梯度下降, 我们得到每个实例 x^t 的更新规则:

$$\Delta m_{ij} = -\eta \frac{\partial E^t}{\partial m_{ij}} = \eta b_i^t (x_j^t - m_{ij}) \quad (12.5)$$

这把最近的中心 (其 $b_i^t=1$) 向输入移动一个因子 η 。其他中心的 $b_l^t (l \neq i)$ 等于 0, 并且不更新 (参见图 12-1)。批处理过程也可以通过将式(12.5)在所有的 t 上求和定义。与任何梯度下降过程一样, 也可以添加一个动量项。为了收敛, η 逐渐减少为 0。但是, 这意味稳定性与可塑性的两难选择 (stability-plasticity dilemma): 如果 η 向 0 递减, 则网络变得稳定, 但是因为更新变得太小, 我们失去了对随时出现的新模式的适应性。如果我们一直保持 η 较大, 则 m_i 可能震荡。

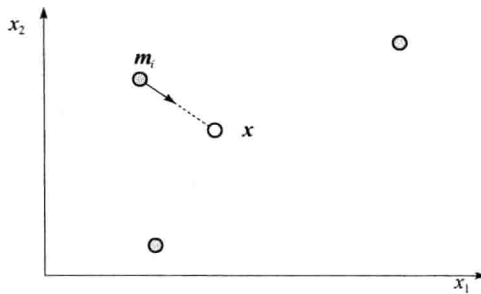


图 12-1 阴影圆是中心, 空心圆是输入实例。 k -均值算法的在线版本沿方向 $(x - m_i)$ 将最近的中心移动一个因子 η

在线 k -均值的伪代码在图 12-2 中。这是图 7-3 的批处理算法的在线版本。

```

初始化  $m_i (i=1, \dots, k)$ , 例如为  $k$  个随机的  $x^t$ 
Repeat
  For 随机次序下所有的  $x^t \in X$ 
     $i \leftarrow \arg \min_j \|x^t - m_j\|$ 
     $m_i \leftarrow m_i + \eta (x^t - m_i)$ 
Until  $m_i$  收敛
  
```

图 12-2 在线 k -均值算法。批处理版本在图 7-3 中

竞争网络可以用单层递归网络实现,如图 12-3 所示。输入层包含输入向量 \mathbf{x} ; 注意没有偏倚单元。输出单元的值是 b_i , 并且它们是感知器:

$$b_i = \mathbf{m}_i^T \mathbf{x} \quad (12.6)$$

然后,我们需要选择最大的 b_i , 并将它设置为 1, 而设置其他 $b_l (l \neq i)$ 为 0。如果我们想用纯粹的神经系统方法做所有的事,即使用并发操作处理单元网络,则最大值的选择可以用横向抑制(lateral inhibition)实现。如图 12-3 所示,每个单元有一个到自身的兴奋的递归连接(即具有正权重)和到其他输出单元的抑制的递归连接(即具有负权重)。使用适当的非线性激励函数和正的及负的递归权重值,这样的网络在某些迭代后收敛于一种状态,其中最大值变成 1,而其余的变成 0 (Grossberg 1980, Feldman 和 Ballard 1982)。

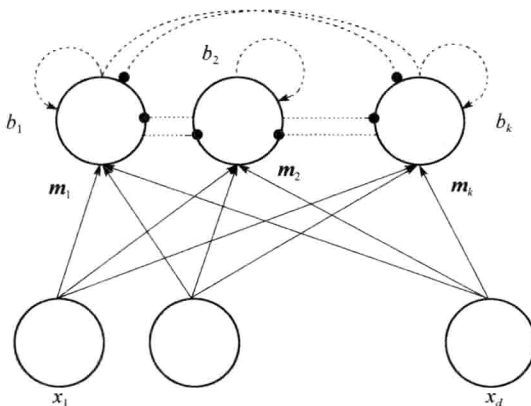


图 12-3 胜者全取竞争神经网络,它是在输出层有递归连接的 k 个感知器的网络。虚线是递归连接,其中带箭头的是兴奋的,而带圆点的是抑制的。输出层的每个单元加强它的值,并试图超过其他输出。在这些递归权重的适当赋值下,最大的抑制了其他所有的。这具有网络效应:其 \mathbf{m}_i 最接近 \mathbf{x} 的一个单元以其 b_i 等于 1 告终,而其他所有的 $b_l (l \neq i)$ 为 0

(12.6) 式中使用的点积是一种相似性度量,并且我们在 5.5 节式(5.26)中看到,如果 \mathbf{m}_i 具有相同的范数,则具有最小欧氏距离 $\|\mathbf{m}_i - \mathbf{x}\|$ 的单元与具有最大点积 $\mathbf{m}_i^T \mathbf{x}$ 的单元相同。

这里和后面,当我们讨论其他竞争方法时,我们使用欧氏距离,但是我们应当记住,使用欧氏距离意味所有输入属性具有相同的方差并且它们是不相关的。如果不是这种情况,则应当反映在距离度量中(即使用马氏距离),或者在使用欧氏距离前,在预处理阶段做适当的规范化(例如用 PCA)。

我们可以将式(12.5)改写为

$$\Delta \mathbf{m}_{ij}^t = \eta b_i^t x_j^t - \eta b_i^t \mathbf{m}_{ij} \quad (12.7)$$

让我们回想一下, \mathbf{m}_{ij} 是从 x_j 到 b_i 的连接权重。正如我们在前一项所看到的,更新

$$\Delta \mathbf{m}_{ij}^t = \eta b_i^t x_j^t \quad (12.8)$$

是 Hebbian 学习(Hebbian learning),它将更新定义为前突触与后突触单元值的乘积。它是作为神经可塑性模型提出的:一个突触变得更重要,如果该连接的前后单元都同时冲动,表明它们是相关的。然而,仅用 Hebbian 学习,权重的增加无界($x_j^t \geq 0$),并且我们需要第二种力量,来减少未更新的权重。一种可能的方法是显式地规范化权重,使得 $\|\mathbf{m}_i\| = 1$; 如果 $\Delta \mathbf{m}_{ij} > 0$, $\Delta \mathbf{m}_{il} = 0 (l \neq i)$, 一旦我们把 \mathbf{m}_i 规范化为单位向量,则 \mathbf{m}_{il} 减少。另一种可能的方

法是引进权衰减项(Oja 1982),而式(12.7)的第二项就可以看作这样的项。Hertz, Krogh 和 Palmer(1991)更详细地讨论了竞争网络和 Hebbian 学习,并且展示这种网络如何学习做 PCA。Mao 和 Jain(1995)讨论了 PCA 和 LDA 的在线算法。

正如我们在第7章所看到的,一个问题是避免死中心,即存在却没有被实际利用的中心。在竞争网络,这对应因为被初始化远离任何输入而从来未能赢得竞争的中心。存在多种方法避免它:

- 1) 我们可以通过随机地选择输入实例来初始化 m_i , 并且确保它们从有数据的地方开始。
- 2) 我们可以使用领导者聚类算法并且逐个添加单元,总是将它们添加在需要它们的地方。一个例子是 ART 模型,我们将在 12.2.2 节讨论。
- 3) 更新时不仅更新最近单元的中心,而且也更新某些其他中心。随着它们被更新,它们也向输入移动,逐渐移向输入空间存在输入的部分,并最终赢得竞争。一个例子是我们将在 12.2.3 节讨论的 SOM。
- 4) 另一种可能引进良心(conscience)机制(Desieno 1988):当前赢得竞争的单元有负罪感并允许其他单元获胜。

12.2.2 自适应共鸣理论

在计算参数之前,应当知道并指定分组数 k 。另一种方法是增量的(incremental),从单个分组开始,并在需要时添加新的分组。作为增量算法的一个例子,我们讨论自适应共鸣理论(adaptive resonance theory, ART)算法(Carpenter 和 Grossberg 1988)。在 ART 中,给定一个输入,所有的输出单元计算它们的值,并且选择与输入最相似的单元。如果使用如式(12.6)中的点积,则它是具有最大值的单元;如果使用欧氏距离,则它是具有最小值的单元。

让我们假设我们使用欧氏距离。如果最小值小于某个称作警戒值(vigilance)的阈值,则像在线 k -均值一样进行更新。如果距离大于警戒值,则增加一个新的输出单元,并且它的中心用该实例初始化。这定义了一个超球,其半径由定义每个单元范围体积的警戒值给定。当我们具有一个不能被任何单元覆盖的输入时,我们就增加一个新单元(参见图 12-4)。

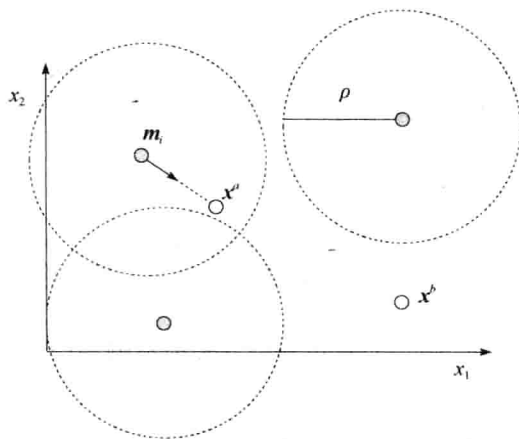


图 12-4 从 x^a 到最近中心的距离小于警戒值 ρ , 中心像在线 k -均值一样进行更新。
然而, x^b 与任何中心都不足够近, 应当在该位置创建一个新的分组

记警戒值为 ρ ，在每次更新我们使用下式：

$$b_i = \| \mathbf{m}_i - \mathbf{x}^t \| = \min_{l=1}^k \| \mathbf{m}_l - \mathbf{x}^t \| \quad (12.9)$$

$$\begin{cases} \mathbf{m}_{k+1} \leftarrow \mathbf{x}^t & \text{如果 } b_i > \rho \\ \Delta \mathbf{m}_i = \eta(\mathbf{x}^t - \mathbf{m}_i) & \text{否则} \end{cases}$$

在距离上设定阈值等价于在每个距离的重构误差上设定阈值，并且如果距离是欧氏距离、误差像式(12.4)那样定义，则这表明每个实例允许的最大重构误差为警戒值的平方。

12.2.3 自组织映射

避免死单元的一种方法是不仅更新获胜者，而且也更新某些其他单元。在 Kohonen (1990, 1995) 提出的自组织映射(self-organizing map, SOM)中，单元下标(即，如 \mathbf{m}_i 中的 i)定义单元的邻域(neighborhood)。当 \mathbf{m}_i 是最近的中心时，除 \mathbf{m}_i 更新之外，还更新它的近邻。例如，如果邻域大小为 2，则 \mathbf{m}_{i-2} , \mathbf{m}_{i-1} , \mathbf{m}_{i+1} , \mathbf{m}_{i+2} 也更新，但是随邻域的加大，使用较小的权重。如果 i 是最近中心的下标，则中心按下式更新

$$\Delta \mathbf{m}_i = \eta e(l, i) (\mathbf{x}^t - \mathbf{m}_i) \quad (12.10)$$

其中 $e(l, i)$ 是邻域函数。当 $l=i$ 时 $e(l, i) = 1$ ，并随 $|l-i|$ 增大而减小。例如，定义它为高斯函数 $\mathcal{N}(i, \sigma)$ ：

$$e(l, i) = \frac{1}{\sqrt{2\pi}\sigma} \exp \left[-\frac{(l-i)^2}{2\sigma^2} \right] \quad (12.11)$$

为了收敛，邻域函数的支集随时间减小，例如 σ 减小，最终只有一个获胜者被更新。

由于邻域单元也向输入移动，我们避免了死单元，因为从它们的近邻朋友得到一点初始帮助之后，稍后的某个时候它们将赢得竞争(参见图 12-5)。

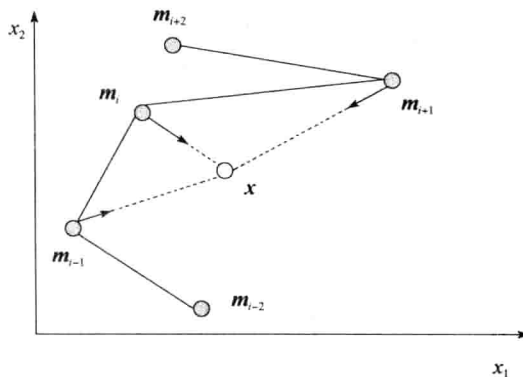


图 12-5 在 SOM 中，不仅最近的单元，而且还有它的近邻(就下标而言)都向输入移动。这里，邻域为 1； \mathbf{m}_i 和它的 1-近邻被更新。注意，这里 \mathbf{m}_{i+1} 远离 \mathbf{m}_i ，但是随着它与 \mathbf{m}_i 一起更新，并且当 \mathbf{m}_{i+1} 是胜者时 \mathbf{m}_i 也被更新，它们最终也成为输入空间的近邻

更新近邻具有如下效果：即使中心被随机初始化，因为它们一起朝着相同的输入移动，因此一旦系统收敛，具有相同邻近下标的单元也将是输入空间中的近邻。

在大部分应用中，单元被组织成二维映射(map)。即每个单元将具有两个下标 $\mathbf{m}_{i,j}$ ，并

且邻域定义在两个维上。如果 $m_{i,j}$ 是最近的中心, 则中心按下式更新

$$\Delta m_{k,l} = \eta e(k,l,i,j) (\mathbf{x}^t - m_{k,l}) \quad (12.12)$$

其中邻域函数现在是二维的。收敛后, 这形成了原 d -维输入空间的二维地形图 (topographical map)。该图包含了空间的高密度部分的许多单元, 而对于没有输入的部分则不显示其中的任何单元。一旦该图收敛, 则原空间靠近的输入被映射到该图中靠近的单元。从这种角度讲, 该图可以解释为做一个非线性形式的多维缩放, 将原来的 \mathbf{x} 空间映射到二维 (i, j) 上。类似地, 如果映射是一维的, 则单元放置在输入空间最大密度的曲线上, 作为主曲线 (principal curve)。

12.3 径向基函数

在隐藏单元使用点积的多层感知器中 (第 11 章), 每一个隐藏单元定义了一个超平面, 并且使用 S 形非线性, 隐藏单元具有 0 和 1 之间的值, 对每个实例关于超平面的位置编码。每个超平面将输入空间一分为二, 并且通常对于给定的输入, 许多隐藏单元都具有非零输出。这称作分布表示 (distributed representation), 因为输入被许多隐藏单元的同时激活重新编码。

另一种可能性是局部表示 (local representation), 其中对于给定的输入, 只有一个或几个单元是活跃的。就像这些局部调整 (locally tuned) 的单元在它们之间划分输入空间并且只对某些输入是有选择的。输入空间的单元中具有非零响应的部分称作接受域 (receptive field)。输入空间则被这样的单元覆盖。

在大脑皮层的多处发现了具有这种响应特征的神经元。例如, 视觉皮层细胞对刺激有选择地响应, 既局部于视网膜的位置, 又局部于视觉方向的角度。这种局部调整的细胞通常排列在大脑皮层图上, 那里像在 SOM 中一样, 细胞对其响应的变量值随它们在图中的位置而变化。

局部性意味有一个距离函数, 度量给定输入 \mathbf{x} 和单元 h 的位置 m_h 的相似度。通常, 该度量取欧氏距离 $\|\mathbf{x} - m_h\|$ 。选取响应函数使得当 $\mathbf{x} = m_h$ 时取最大值, 并且随着它们的相似性减小而减少。通常, 我们使用高斯函数 (参见图 12-6):

$$p_h^t = \exp \left[-\frac{\|\mathbf{x}^t - m_h\|^2}{2s_h^2} \right] \quad (12.13)$$

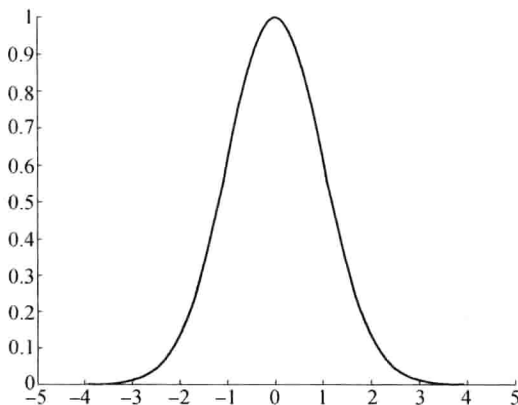


图 12-6 用于径向基函数网络的钟形函数的一维形式。这个函数 $m=0$, $s=1$ 。它像高斯函数但不是密度函数; 其积分不等于 1。在 $(m-3s, m+3s)$ 中它非零, 但是更保守的区间是 $(m-2s, m+2s)$

严格地说,这不是高斯密度,但是我们还是使用了相同的名字。 m_j 和 s_j 分别表示局部单元 j 的中心和展宽,这样定义了一个径向对称的基函数。以使用更复杂的模型为代价,我们可以使用椭球,不同的维上具有不同的展宽,甚至使用马氏距离,允许相关的输入。

使用局部基函数的基本思想是,在输入数据中存在一些实例分组或簇,而对每个簇,我们定义一个基函数 p_h^i ,当实例 x^i 属于簇 h 时它不为零。我们可以使用 12.2 节讨论的任意在线竞争方法来找出中心 m_h 。有一种简单、有效的启发式方法来找出展宽:一旦我们求出中心,我们就可以找出簇中的最远实例,并令 s_h 为它到簇中心距离的一半。我们本来也可以使用三分之一,但是我们宁愿保守一点。我们还可以使用统计聚类方法找出簇参数。例如,在高斯混合分布上使用第 7 章讨论的 EM 方法,找出簇参数均值、方差(和协方差)。

$p_h^i (h=1, \dots, H)$ 定义了一个新的 H -维空间,并且形成 x^i 的新表示。我们也可以使用 b_h^i 式(12.2)对输入编码,但是 b_h^i 为 0/1; p_h^i 具有附加的优点,它用 $(0, 1)$ 中的值对点到其中心的距离编码。该值衰减到 0 的速度依赖于 s_h 。图 12-7 给出了一个例子,并且将这种局部表示与多层感知器使用的分布表示进行比较。由于高斯分布都是局部的,因此与使用分布表示相比,通常我们需要更多的局部单元,当输入是高维的时尤其如此。

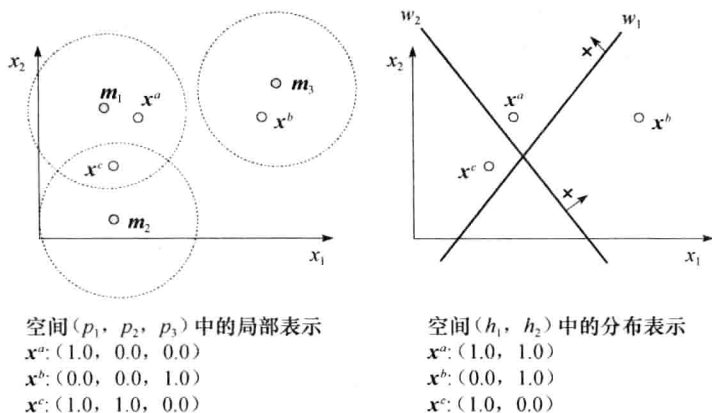


图 12-7 局部表示与分布表示之间的差别。值是硬的 0/1 值。我们可以使用 $(0, 1)$ 之间的软值得到更多信息编码。在局部表示中,用高斯 RBF 做,它使用到中心 m_i 的距离;而在分布表示中,使用 S 形函数做,它使用到超平面 w_i 的距离

在监督学习的情况下,我们可以使用这种新的局部表示作为输入。如果我们使用感知器,则我们有

$$y^t = \sum_{h=1}^H w_h p_h^t + w_0 \quad (12.14)$$

其中 H 是基函数的个数。这种结构称作径向基函数(radial basis function, RBF)网络(Broomhead 和 Lowe 1988; Moody 和 Darken 1989)。通常,人们不使用多于一个高斯单元层的 RBF 网络。 H 是复杂度参数,像多层感知器的隐藏单元数一样。之前,当它对应非监督学习中的中心数时,我们用 k 表示它。

这里,我们看到使用 p_h 而不是使用 b_h 的优点。由于 b_h 是 0/1,如果式(12.14)使用 b_h 而不是 p_h ,则它将给出分段常量近似,在单元区域的边界不连续。 p_h 值是软的并导致光滑的近似,从一个区域到另一个时取加权平均。我们可以容易地看到这种网络是一种普适近似,因为

给定足够多的单元，它可以以期望的精度逼近任意函数。对于我们期望的精度，我们可以形成一个输入空间的网格，实际上对每个网格定义一个单元，并设置它的权重 w_h 为期望的输出值。

这种结构与非参数估计(例如我们在第8章所看到的 Parzen 窗口)很相似，并且 p_h 可以看作核函数。不同之处在于我们没有整个训练集上的核函数，而是使用聚类方法将它们分组，使用更少的核函数。单元数 H 是复杂度参数，在简洁性和准确性之间平衡。使用更多单元，我们就更好地逼近训练数据，但是我们得到更复杂的模型并有过分拟合的风险；太少可能拟合不足。最佳值仍然用交叉验证确定。

一旦给定和固定 m_h 和 s_h ， p_h 也是固定的。然后可以容易地批处理或在线地训练 w_h 。对于回归，这是一个线性回归模型(p_h 作为输入)，并且 w_h 可以解析地求解，而不需要迭代(参见4.6节)。对于分类，我们需要借助于迭代过程。我们在第10章讨论过这些学习方法，此处不再赘述。

这里，我们要做的是个两阶段过程：我们使用非监督方法确定中心，然后在其上构建监督层。这称作混合学习(hybrid learning)。我们还可以用监督的方式学习所有的参数，包括 m_h 和 s_h 。式(12.13)的径向基函数是可微的，并且我们可以后向传播，就像我们在多层感知器后向传播，更新第一层的权重一样。该结构类似于多层感知器，以 p_h 为隐藏单元， m_h 和 s_h 作为第一层的参数，高斯函数作为隐藏层的激励函数，而 w_h 作为第二个隐藏层的权重(见图12-8)。

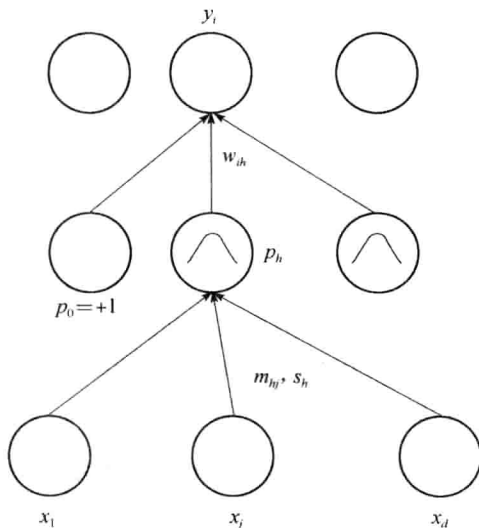


图 12-8 RBF 网络，其中 p_h 是使用钟形激励函数的隐藏单元， m_h, s_h 是第一层的参数，而 w_h 是第二层的权重

但是，在我们讨论这些之前，我们应当注意到训练两层网络很慢。混合模型一次训练一层，因而比较快。另一种技术，称作锚(anchor)方法，将中心设置为从训练集随机选取的模式，而不进一步更新。如果有许多单元，这足以满足需要。

另一方面，精度通常没有使用完全监督方法高。考虑输入是均匀分布的情况。 k -均值聚类均匀地安放单元。如果函数在一小部分空间稍有变化，则更好的想法是将更多单元安放在函数变化快的地方，以便使误差尽可能小；这正是完全监督方法所要做的。

让我们讨论如何在完全监督方式下训练所有参数。方法与用于多层感知器的后向传播一样。让我们考虑具有多个输出的回归。批处理的误差为

$$E(\{\mathbf{m}_h, s_h, w_{ih}\}_{i,h} | \mathcal{X}) = \frac{1}{2} \sum_i \sum_h (r'_i - y'_i)^2 \quad (12.15)$$

其中

$$y'_i = \sum_{h=1}^n w_{ih} p'_h + w_{i0} \quad (12.16)$$

使用梯度下降，我们得到第二层权重的如下更新规则：

$$\Delta w_{ih} = \eta \sum_i (r'_i - y'_i) p'_h \quad (12.17)$$

这是通常的感知器更新规则，其中 p_h 作为输入。通常， p_h 交叠不多，并且在每次迭代中，只有少量的 p_h 非零，并且只有它们的 w_h 被更新。这就是为什么 RBF 网络学习非常快，并且比使用分布表示的多层感知器快的原因。

类似地，我们可以用后向传播（链规则）得到中心和展宽的更新方程：

$$\Delta m_{hj} = \eta \sum_i \left[\sum_h (r'_i - y'_i) w_{ih} \right] p'_h \frac{(x'_j - m_{hj})}{s_h^2} \quad (12.18)$$

$$\Delta s_h = \eta \sum_i \left[\sum_i (r'_i - y'_i) w_{ih} \right] p'_h \frac{\|\mathbf{x}' - \mathbf{m}_h\|^2}{s_h^3} \quad (12.19)$$

让我们比较式(12.18)和式(12.5)：首先，这里我们使用 p_h 而不是 b_h ，这意味不仅最近的单元而且所有的单元都根据它们的中心和展宽而被更新。其次，这里更新是监督的，并且包含后向传播的误差项。更新不仅依赖于输入，而且依赖于最终的误差 $(r'_i - y'_i)$ 、单元对输出的影响 w_{ih} 、单元的活性 p_h 和输入 $(\mathbf{x} - \mathbf{m}_i)$ 。

对于分类，我们有

$$y'_i = \frac{\exp \left[\sum_h w_{ih} p'_h + w_{i0} \right]}{\sum_k \exp \left[\sum_h w_{kh} p'_h + w_{k0} \right]} \quad (12.20)$$

而互熵误差为

$$E(\{\mathbf{m}_h, s_h, w_{ih}\}_{i,h} | \mathcal{X}) = - \sum_i \sum_h r'_i \log y'_i \quad (12.21)$$

使用梯度下降，可以类似地导出更新规则（习题3）。

让我们再考虑式(12.14)：对于任意输入，如果 p_h 非零，则它对输出的贡献为 w_h 。它的贡献是常量拟合，由 w_h 给定。通常，高斯函数交叠不多，并且它们之中的一或两个具有非零的 p_h 值。在任何情况下，只有少数单元对输出有贡献。 w_0 是偏移常量，加到活跃（非零）单元的加权和上。我们还看到如果所有的 p_h 均为0，则 $y = w_0$ 。这样，我们可以把 w_0 看作是 y 的缺省值：如果没有高斯单元是活跃的，则输出由该值给定。因此，有可能使得该“默认模型”具有更强的能力。例如，我们可以令

$$y' = \sum_{h=1}^n w_h p'_h + \mathbf{v}^T \mathbf{x}' + v_0 \quad (12.22)$$

在这种情况下，默认模型是线性的： $\mathbf{v}^T \mathbf{x}' + v_0$ 。当它们非零时，高斯模型好像是“例外”

并修改输出, 补偿期望输出与默认输出之差。这种模型可以以监督方式训练, 而默认模型与 w_h 一起训练(习题 4)我们将在 17.11 节讨论类似的模型级联(cascading), 其中我们将看到两个学习器的组合, 一个是一般规则, 另一个由一组异常组成。

12.4 结合基于规则的知识

如果我们能够结合先验知识(prior knowledge)来初始化, 则任何学习系统的训练都可以更简单。例如, 先验知识可以以一组规则的形式提供, 指定该模型(例如 RBF 网络)必须学习的输入/输出映射。这种情况在业界和医学应用中经常出现, 在应用中规则可以由专家提供。类似地, 一旦网络被训练, 就可以从中提取规则, 使得问题的解更容易理解。

包含先验知识还有其他优点。如果需要将网络外推到输入空间中从未见到训练数据的区域, 可以依赖这种先验知识。此外, 在许多控制应用中, 需要网络一开始就做出合理的预测。在它看到足够多的训练数据之前, 必须主要依赖这种先验知识。

在许多应用中, 我们通常被告知一些基本规则, 开始需要遵循, 而后通过经验来精炼和改变。我们关于问题的初始知识越好, 我们得到好性能就越快, 并且需要的训练就越少。

使用 RBF 网络, 这种包含先验知识或提取学习的知识很容易做, 因为单元是局部的。这使得规则提取(rule extraction)更容易(Tresp、Hollatz 和 Ahmad 1997)。一个例子是

$$\text{IF}((x_1 \approx a) \text{ AND } (x_2 \approx b)) \text{ OR } (x_3 \approx c) \text{ THEN } y = 0.1 \quad (12.23)$$

其中 $x_1 \approx a$ 意指“ x_1 约等于 a ”。在 RBF 框架中, 这个规则被两个高斯单元编码为

$$p_1 = \exp\left[-\frac{(x_1 - a)^2}{2s_1^2}\right] \cdot \exp\left[-\frac{(x_2 - b)^2}{2s_2^2}\right], \quad \text{其中 } w_1 = 0.1$$

$$p_2 = \exp\left[-\frac{(x_3 - c)^2}{2s_3^2}\right], \quad \text{其中 } w_2 = 0.1$$

“约等于”被一个高斯函数建模, 这里中心是理想值, 展宽表示理想值周围允许的差。合取是两个一元高斯函数的积, 是二元高斯函数。于是, 第一个乘积项可以被二维(即 $\mathbf{x} = [x_1, x_2]$)高斯函数处理, 其中心在 (a, b) , 而在两个维上的展宽由 s_1 和 s_2 给定。析取被两个单独的高斯函数建模, 每个处理一个析取项。

给定标记的训练数据, 使用较小的 η 值, 这样构造的 RBF 网络的参数在初始构造后可以微调。

这种表示方法与模糊逻辑方法有关, 那里式(12.23)称作模糊规则(fuzzy rule)。检查近似相等的高斯基函数对应模糊隶属关系函数(fuzzy membership function)(Berthold 1999; Cherkassky 和 Mulier 1998)。

12.5 规范化基函数

在式(12.14)中, 对于一个输入, 可能所有的 p_h 都为 0。在某些应用中, 我们可能希望有一个规范化步骤, 确保局部单元值的和为 1, 从而确保对于任何输入, 至少存在一个非零单元:

$$g_h^i = \frac{p_h^i}{\sum_{l=1}^H p_l^i} = \frac{\exp[-\|\mathbf{x}^i - \mathbf{m}_h\|^2 / 2s_h^2]}{\sum_l \exp[-\|\mathbf{x}^i - \mathbf{m}_l\|^2 / 2s_l^2]} \quad (12.24) \quad \boxed{295}$$

图 12-9 给出了一个例子。取 p_h 为 $p(\mathbf{x} | h)$, g_h 对应 \mathbf{x} 属于单元 h 的后验概率 $p(h | \mathbf{x})$ 。这就像单元在它们之间划分输入空间。我们可以想象 g_h 本身是分类器, 为给定的输入选择响应单元。这种分类基于距离来做, 就像在有参高斯分类器中那样(第 5 章)。

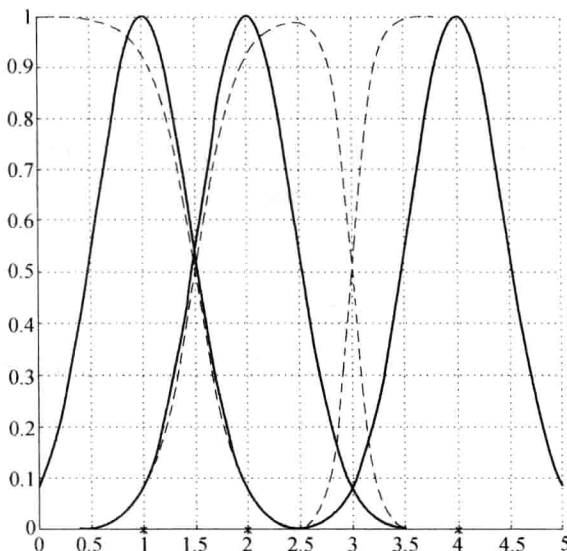


图 12-9 规范化前 (—) 和规范化后 (---) 的三个高斯分布, 其中心用“*”标记。注意一个单元的非零区域还依赖其他单元的位置。如果展宽较小, 则规范化实现较硬的划分; 使用较大的展宽, 单元重叠更多

输出是加权和

$$y_i^t = \sum_{h=1}^H w_{ih} g_h^i \quad (12.25)$$

这里不需要偏倚项, 因为对于每个 \mathbf{x} , 至少有一个非零的 g_h 。使用 g_h 而不是 p_h 并不引进附加的参数; 它只是将单元联系在一起: p_h 仅依赖于 \mathbf{m}_h 和 s_h , 但是由于规范化, g_h 依赖于所有单元的中心和展宽。

对于回归, 使用梯度下降, 我们有如下更新规则:

$$\Delta w_{ih} = \eta \sum_t (r_i^t - y_i^t) g_h^i \quad (12.26)$$

$$\Delta m_{hj} = \eta \sum_t \sum_i (r_i^t - y_i^t) (w_{ih} - y_i^t) g_h^i \frac{(x_j^t - m_{hj})}{s_h^2} \quad (12.27)$$

可以类似地导出 s_h 的更新规则和用于分类的规则。让我们把这些规则与具有非规范化高斯分布的 RBF 的规则式(12.17)进行比较。这里, 我们使用 g_h 而不是 p_h , 这使得单元的更新不仅依赖于它自己的参数, 而且也依赖于其他单元的中心和展宽。比较式(12.27)和式(12.18), 我们看到我们有 $(w_{ih} - y_i^t)$ 而不是 w_{ih} , 这展示了规范化在输出上的作用。“负责任”的单元希望降低它的输出 w_{ih} 和最终输出 y_i^t 之间的差, 正比于它的责任 g_h 。

12.6 竞争的基函数

正如我们迄今为止所看到的, 在 RBF 网络中, 最终的输出被局部单元贡献的加权和所确定。尽管单元是局部的, 但是重要的是最终的加权和, 并且我们希望使它与要求的输出尽可能接近。例如对于回归, 我们最小化式(12.15), 这基于概率模型

$$p(\mathbf{r}^t | \mathbf{x}^t) = \prod_i \frac{1}{\sqrt{2\pi}\sigma} \exp \left[-\frac{(r_i^t - y_i^t)^2}{2\sigma^2} \right] \quad (12.28)$$

其中 y_i^t 由式(12.16) (非规范化的) 或式(12.25) (规范化的) 给出。在两种情况下, 我们都可以将模型看作协同 (cooperative) 模型, 因为单元协同操作, 以便产生最终的输出 y_i^t 。现在, 我们讨论使用竞争基函数 (competitive basis function) 的方法, 其中我们假定输出取自混合模型

$$p(\mathbf{r}^t | \mathbf{x}^t) = \sum_{h=1}^H p(h | \mathbf{x}^t) p(\mathbf{r}^t | h, \mathbf{x}^t) \quad (12.29)$$

$p(h | \mathbf{x}^t)$ 是混合比例, $p(\mathbf{r}^t | h, \mathbf{x}^t)$ 是产生输出的混合分支, 如果该分支被选择的话。注意这两项都依赖于输入 \mathbf{x} 。

混合比例为

$$p(h | \mathbf{x}) = \frac{p(\mathbf{x} | h)p(h)}{\sum_l p(\mathbf{x} | l)p(l)} \quad (12.30)$$

$$g_h^t = \frac{a_h \exp[-\|\mathbf{x}^t - \mathbf{m}_h\|^2 / 2s_h^2]}{\sum_l a_l \exp[-\|\mathbf{x}^t - \mathbf{m}_l\|^2 / 2s_l^2]} \quad (12.31)$$

一般, 我们假定 a_h 相等并忽略它们。让我们先考虑回归, 其中分支是高斯的。在式(12.28), 噪声被加到加权和上; 这里一个分支被选中, 并且噪声加到它的输出 y_{ih}^t 上。

使用式(12.29)的混合模型, 对数似然是

$$\mathcal{L}(\{\mathbf{m}_h, s_h, w_{ih}\}_{i,h} | \mathcal{X}) = \sum_i \log \sum_h g_h^t \exp \left[-\frac{1}{2} \sum_i (r_i^t - y_{ih}^t)^2 \right] \quad (12.32)$$

其中 $y_{ih}^t = w_{ih}$ 是由分支 h 对输出 i 做的常量拟合。严格地说, 它不依赖于 \mathbf{x} 。(在 12.8.2 节, 我们讨论竞争的混合专家模型, 其中局部拟合是 \mathbf{x} 的线性函数。)我们看到如果 g_h^t 为 1, 则它对产生正确的输出负责, 并且需要最小化它的预测误差的平方和 $\sum_i (r_i^t - y_{ih}^t)^2$ 。

使用梯度上升最大化该对数似然, 我们得到

$$\Delta w_{ih} = \eta \sum_t (r_i^t - y_{ih}^t) f_h^t \quad (12.33)$$

其中

$$f_h^t = \frac{g_h^t \exp \left[-\frac{1}{2} \sum_i (r_i^t - y_{ih}^t)^2 \right]}{\sum_l g_l^t \exp \left[-\frac{1}{2} \sum_i (r_i^t - y_{il}^t)^2 \right]} \quad (12.34)$$

$$p(h | \mathbf{r}, \mathbf{x}) = \frac{p(h | \mathbf{x})p(\mathbf{r} | h, \mathbf{x})}{\sum_l p(l | \mathbf{x})p(\mathbf{r} | l, \mathbf{x})} \quad (12.35)$$

$g_h^t \equiv p(h | \mathbf{x}^t)$ 是给定输入、单元 h 的后验概率，并且它依赖于所有单元的中心和展宽。 $f_h^t \equiv p(h | \mathbf{r}, \mathbf{x}^t)$ 是给定输入和期望的输出、单元 h 的后验概率，也在选择负责单元时考虑误差。

类似地，我们可以推导更新中心的规则：

$$\Delta m_{hj} = \eta \sum_t (f_h^t - g_h^t) \frac{(x_j^t - m_{hj})}{s_h^2} \quad (12.36)$$

f_h 是单元 h 的后验概率，也考虑要求的输出，而 g_h 是仅使用输入空间信息的后验概率。它们的差是中心的误差项。可以类似地导出 Δs_h 。在协同情况下，并不强求单元是局部的。为了降低误差，均值和展宽都可以取任意值；有时，甚至可以增加和展平展宽。然而，在竞争情况下，为了提高似然，单元必须是局部的，它们之间更加分离，并具有更小的展宽。

对于分类，每个分支本身是多项式。于是，对数似然为

$$\mathcal{L}(\{\mathbf{m}_h, s_h, w_{ih}\}_{i,h} | \mathcal{X}) = \sum_t \log \sum_h g_h^t \prod_i (y_{ih}^t)^{r_i} \quad (12.37)$$

$$= \sum_t \log \sum_h g_h^t \exp \left[\sum_i r_i^t \log y_{ih}^t \right] \quad (12.38)$$

其中

$$y_{ih}^t = \frac{\exp w_{ih}}{\sum_k \exp w_{kh}} \quad (12.39)$$

可以使用梯度上升导出 w_{ih} 、 \mathbf{m}_h 和 s_h 的更新规则，这包括

$$f_h^t = \frac{g_h^t \exp \left[\sum_i r_i^t \log y_{ih}^t \right]}{\sum_l g_l^t \exp \left[\sum_i r_i^t \log y_{il}^t \right]} \quad (12.40)$$

在第7章，我们讨论了用混合高斯模型拟合数据的 EM 算法。也可以将 EM 推广到监督学习。实际上，计算 f_h^t 对应 E 步。 $f_h^t \equiv p(\mathbf{r} | h, \mathbf{x}^t)$ 取代了 $p(h | \mathbf{x}^t)$ ，后者是应用处于非监督状态时我们在第7章的 E 步所使用的。对于回归，在 M 步我们用下式更新参数

$$\mathbf{m}_h = \frac{\sum_t f_h^t \mathbf{x}^t}{\sum_t f_h^t} \quad (12.41)$$

$$s_h = \frac{\sum_t f_h^t (\mathbf{x}^t - \mathbf{m}_h)(\mathbf{x}^t - \mathbf{m}_h)^T}{\sum_t f_h^t} \quad (12.42)$$

$$w_{ih} = \frac{\sum_t f_h^t r_i^t}{\sum_t f_h^t} \quad (12.43)$$

我们看到 w_{ih} 是加权平均，其中权重是给定输入和期望的输出、单元的后验概率。对于分类，M 步没有解析解并且需要借助于迭代过程，例如梯度上升 (Jordan 和 Jacobs 1994)。

12.7 学习向量量化

假设对每个类我们有 H 个单元, 已经被这些类标记。这些单元已经被它们的类中的实例随机初始化。在每次迭代中, 我们发现单元 m_i 在欧式距离下最接近输入实例, 并使用如下更新规则:

$$\begin{cases} \Delta m_i = \eta(x^t - m_i) & \text{如果 } x^t \text{ 和 } m_i \text{ 具有相同的类标号} \\ \Delta m_i = -\eta(x^t - m_i) & \text{否则} \end{cases} \quad (12.44)$$

如果最近的中心具有正确的标号, 则它将移向输入, 以便更好地代表它。如果它属于错误的类, 则它远离输入, 我们期望的情形是, 如果它移得足够远, 则在未来的迭代中正确的类将是最近的。这称作学习向量量化(learning vector quantization, LVQ)模型, 由 Kohonen 提出(1990, 1995)。

LVQ 更新方程类似于式(12.36), 其中中心移动的方向依赖于两个值的差: 获胜单元基于输入距离的预测和获胜者基于要求的输出。

12.8 混合专家模型

在 RBF 中, 对应每个局部小片, 我们给出一个常量拟合。在对于任意输入, 我们有一个 g_h 为 1, 而其余为 0 的情况下, 我们得到一个分段常量近似, 其中对于输出 i , 小片 h 的局部拟合由 w_{ih} 给出。从泰勒展开式我们知道在每个点, 函数可以写成

$$f(x) = f(a) + (x - a)f'(a) + \cdots \quad (12.45)$$

这样, 如果 x 足够接近 a 并且 $f'(a)$ 接近 0, 即如果 $f(x)$ 在 a 附近是平坦的, 则常量近似很好。如果不是这种情况, 则需要将空间划分成大量小片。当输入维度很高时, 由于维灾难的缘故, 这将是一个特别严重的问题。

一种可供选择的方法是考虑泰勒展开式的下一项(即线性项), 使用分段线性近似(piecewise linear approximation)。这就是混合专家模型(mixture of experts)所做的(Jacobs 等, 1991)。我们令

$$y_i^t = \sum_{h=1}^H w_{ih} g_h^t \quad (12.46)$$

它与式(12.25)一样, 但是这里小片 h 对输出 i 的贡献 w_{ih} 不是常量, 而是输入的线性函数:

$$w_{ih}^t = v_{ih}^T x^t \quad (12.47)$$

v_{ih} 是参数向量, 定义线性函数, 并包含一个偏倚项, 使得混合专家模型是 RBF 网络的推广。单元活性可以取规范化的 RBF:

$$g_h^t = \frac{\exp[-\|x^t - m_h\|^2 / 2s_h^2]}{\sum_l \exp[-\|x^t - m_l\|^2 / 2s_l^2]} \quad (12.48)$$

除第二层权重不是常量而是线性模型的输出外, 这可以看作 RBF 网络(参见图 12-10)。Jacobs 等(1991)用另一种方法来看它: 他们将 w_h 看作线性模型, 每个都取输入, 并称它们为专家。 g_h 被看作一个门网络(gating network)的输出。门网络就像其输出之和为 1 的分类器一样, 将输入指派给一个专家(参见图 12-11)。

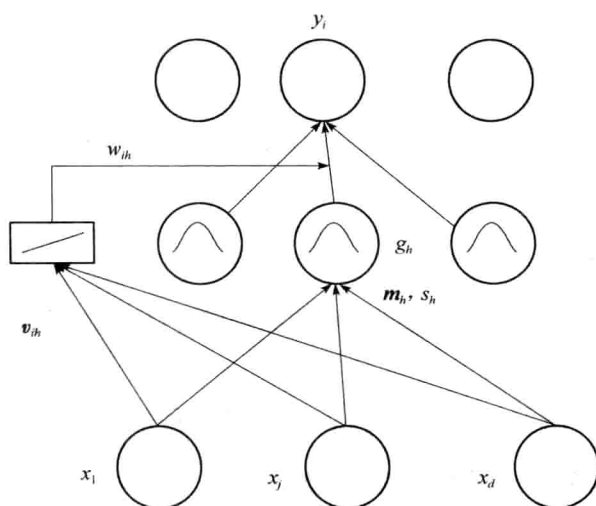


图 12-10 混合专家模型可以看作 RBF 网络，其中第二层的权重是线性模型的输出。为清晰起见，只显示了一个线性模型

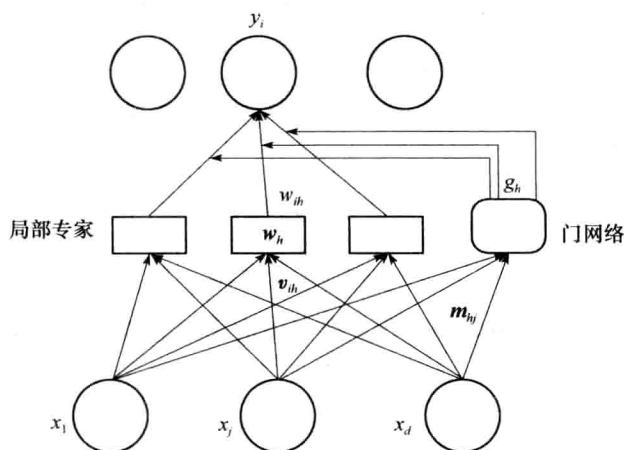


图 12-11 混合专家模型可以看作组合多种模型的模型。 w_h 是模型，而门网络是另一种确定每个模型的权重的模型，如 g_h 所示。从这个角度来看，专家和门网络都不限于是线性的

用这种方式看待门网络，任何分类器都可以用在门网络之中。当 \mathbf{x} 是高维的时，使用局部高斯单元可能需要大量专家，而 Jacobs 等(1991)提议取

$$g_h^t = \frac{\exp[\mathbf{m}_h^T \mathbf{x}^t]}{\sum_l \exp[\mathbf{m}_l^T \mathbf{x}^t]} \quad (12.49)$$

这是一个线性分类器。注意， \mathbf{m}_h 不再是中心，而是超平面，并因此包含偏倚值。门网络实现了分类，它将输入区域线性地划分成专家 h 负责的区域和其他专家负责的区域。正如我们将在第 17 章再次看到的，混合专家模型是一种组合多个模型的一般结构；专家和门网络都可以是非线性的，例如，包含多层感知器而不是线性感知器(见习题 6)。

Bottou 和 Vapnik(1992)提出了一种类似于混合专家模型并进行线性光滑的结构。在他们

的方法中,最初时并不做训练。当给定一个检验样本时,选择一个接近检验实例的数据子集(与 k -最近邻一样,但使用更大的 k),使用这些局部数据训练一个简单模型(如线性分类器),对实例做出预测,然后丢弃该模型。对于下一个实例,创建一个新模型,如此下去。在手写数字识别应用中,这种模型比多层感知器、 k -最近邻和 Parzen 窗口具有更小的误差。缺点是需要现场对每个检验实例训练一个新模型。

12.8.1 协同专家模型

在协同情况下, y_i^t 由式(12.46)给定,而我们希望使它与要求的输出 r_i^t 尽可能接近。对于回归,误差函数是

$$E(\{\mathbf{m}_h, s_h, w_{ih}\}_{i,h} | \mathcal{X}) = \frac{1}{2} \sum_i \sum_t (r_i^t - y_i^t)^2 \quad (12.50)$$

使用梯度下降,则第二层(专家)权重参数更新为

$$\Delta \mathbf{v}_{ih} = \eta \sum_t (r_i^t - y_i^t) \mathbf{g}_h^t \mathbf{x}^t \quad (12.51)$$

与式(12.26)比较,我们看到唯一的区别是,新的更新是输入的函数。

如果我们使用软最大门(参见式(12.49)),则使用梯度下降我们有超平面的如下更新规则:

$$\Delta m_{hj} = \eta \sum_i \sum_t (r_i^t - y_i^t) (w_{ih}^t - y_i^t) \mathbf{g}_h^t \mathbf{x}_j^t \quad (12.52)$$

如果我们使用径向门(参见式(12.48)),则只有最后一项 $\partial p_h / \partial m_{hj}$ 不同。

对于分类,我们有

$$y_i = \frac{\exp \left[\sum_h w_{ih} \mathbf{g}_h^t \right]}{\sum_k \exp \left[\sum_h w_{kh} \mathbf{g}_h^t \right]} \quad (12.53)$$

其中 $w_{ih} = \mathbf{v}_{ih}^T \mathbf{x}$,并且使用梯度下降最小化互熵可以推导出更新规则(习题7)。

12.8.2 竞争专家模型

正如竞争的 RBF,我们有

$$\mathcal{L}(\{\mathbf{m}_h, s_h, w_{ih}\}_{i,h} | \mathcal{X}) = \sum_t \log \sum_h \mathbf{g}_h^t \exp \left[-\frac{1}{2} \sum_i (r_i^t - y_{ih}^t)^2 \right] \quad (12.54)$$

其中 $y_{ih}^t = w_{ih}^t = \mathbf{v}_{ih}^T \mathbf{x}^t$ 。使用梯度上升,我们得到

$$\Delta \mathbf{v}_{ih} = \eta \sum_t (r_i^t - y_{ih}^t) \mathbf{f}_h^t \mathbf{x}^t \quad (12.55)$$

$$\Delta \mathbf{m}_h = \eta \sum_t (f_h^t - g_h^t) \mathbf{x}^t \quad (12.56)$$

假定软最大门如式(12.49)。

对于分类,我们有

$$\mathcal{L}(\{\mathbf{m}_h, s_h, w_{ih}\}_{i,h} | \mathcal{X}) = \sum_t \log \sum_h \mathbf{g}_h^t \prod_i (y_{ih}^t)^{r_i^t} \quad (12.57)$$

$$= \sum_t \log \sum_h \mathbf{g}_h^t \exp \left[\sum_i r_i^t \log y_{ih}^t \right] \quad (12.58)$$

其中

$$y_{ih}^t = \frac{\exp w_{ih}^t}{\sum_k \exp w_{kh}^t} = \frac{\exp[v_{ih} \mathbf{x}^t]}{\sum_k \exp[v_{kh} \mathbf{x}^t]} \quad (12.59)$$

Jordan 和 Jacobs(1994) 将 EM 推广到具有局部线性模型的竞争情况。Alpaydm 和 Jordan(1996) 比较了用于分类人物的协同和竞争模型, 发现协同模型一般更精确, 但是竞争版本学习更快。这是因为在协同情况下, 重叠更多并且实现了比较光滑的近似, 因此更适合回归问题。竞争模型做比较硬的划分; 通常, 对于一个输入, 只有一个专家是活跃的, 因此学习更快。

12.9 层次混合专家模型

在图 12-11 中, 我们看到一组专家和一个选择一个专家作为输入的函数的门网络。在层次混合专家模型(hierarchical mixture of expert)中, 我们以递归的方式用一个完整的混合专家系统取代每个专家(Jordan 和 Jacobs 1994)。这种结构可以看作是一棵决策树(第 9 章), 其中门网络可以看作决策结点。当门网络是线性的时, 这很像第 9.6 节讨论的线性多元决策树。区别是门网络并不做硬决策, 而是取来自子女贡献的加权和。叶结点是线性模型, 并且它们的决策被取平均值并在树中向上传播。树根给出最终的输出, 输出是所有树叶决策的加权和。这是一棵软决策树(soft decision tree), 与我们以前看到的决策树不同, 其中只取一条从树根到树叶的路径。

一旦结构选定, 即选定深度、专家和门模型, 整棵树就可以从标记的样本中学习。Jordan 和 Jacobs(1994) 为这样的结构推导出了梯度下降和 EM 学习规则。

12.10 注释

RBF 网络可以看作神经网络, 由简单处理单元的网络实现。它不同于多层感知器, 第一层和第二层实现了不同的函数。Omohundro(1987) 讨论了如何用神经网络实现局部模型, 并且还提出了相关局部单元快速局部化的层次数据结构。Specht(1991) 证明 Parzen 窗口可以作为神经网络实现。

Platt(1991) 提出了 RBF 的增量版本, 那里新单元可以在必要时添加。类似地, Fritzke(1995) 提出了 SOM 的增长版本。

Lee(1991) 在手写数字识别应用上比较了 k -最近邻、多层感知器和 RBF 网络, 并且结论是三种方法都具有小误差率。RBF 网络学习比多层感知器上的后向传播快, 但是使用更多参数。就分类速度和存储需求而言, 这两种方法都优于 k -NN。在实际应用中, 像时间、存储量和计算复杂性等实际限制可能比误差率的些许差别更重要。

Kohonen 的 SOM(1990, 1995) 是最流行的神经网络方法之一, 已经用于各种各样的应用, 包括探测式数据分析和作为监督学习之前的预处理步骤。一个有趣和成功应用的例子是旅行商人的问题(Angeniol、Vauboiss 和 Le Texier 1988)。正如 k 均值聚类与高斯混合中 EM(参见第 7 章) 之间的区别一样, 生成地形图映射(generative topographic mapping, GTM)

(Bishop, Svebsén 和 Williams 1998) 是 SOM 的概率版本, 它使用其均值落在二维流形上的 (关于低维上的拓扑序) 高斯混合优化数据的对数似然。

在 RBF 网络中, 一旦中心和展览固定 (例如, 像锚方法中那样, 通过随机选择训练实例中的一个子集作为中心), 训练第二层是一个线性模型。这个模型等价于取高斯核的支持向量机。其中在学习期间, 选择最佳的实例子集 (称作支持向量), 我们将在第 13 章讨论。高斯过程 (参见第 14 章) 也类似, 其中我们由存放的训练实例插值。

12.11 习题

1. 给出一个实现 XOR 的 RBF 网络。
2. 写下 RBF 网络, 它使用椭圆单元, 而不是像式 (12.13) 中那样使用径向单元。
3. 为分类的 RBF 网络推导更新方程 (参见式 (12.20) 和式 (12.21))。
4. 展示如何训练式 (12.22) 给定的系统。
5. 比较混合专家结构和 RBF 网络的参数个数。
6. 给出混合专家结构的公式, 其中专家和门网络都是多层感知器。为回归和预测推导更新方程。
7. 为分类推导协同混合专家模型的更新方程。
8. 为分类推导竞争混合专家模型的更新方程。
9. 给出具有两个层次混合专家结构的公式。使用梯度下降, 为回归和分类推导更新方程。
10. 在混合专家模型中, 由于不同的专家专攻输入空间的不同部分, 因此他们需要关注不同的输入。讨论如何在专家中局部地进行维归约。

12.12 参考文献

- Alpaydın, E., and M. I. Jordan. 1996. "Local Linear Perceptrons for Classification." *IEEE Transactions on Neural Networks* 7: 788-792.
- Angeniol, B., G. Vauboiss, and Y. Le Texier. 1988. "Self Organizing Feature Maps and the Travelling Salesman Problem." *Neural Networks* 1: 289-293.
- Berthold, M. 1999. "Fuzzy Logic." In *Intelligent Data Analysis: An Introduction*, ed. M. Berthold and D. J. Hand, 269-298. Berlin: Springer.
- Bishop, C. M., M. Svensén, and C. K. I. Williams. 1998. "GTM: The Generative Topographic Mapping." *Neural Computation* 10: 215-234.
- Bottou, L., and V. Vapnik. 1992. "Local Learning Algorithms." *Neural Computation* 4: 888-900.
- Broomhead, D. S., and D. Lowe. 1988. "Multivariable Functional Interpolation and Adaptive Networks." *Complex Systems* 2: 321-355.
- Carpenter, G. A., and S. Grossberg. 1988. "The ART of Adaptive Pattern Recognition by a Self-Organizing Neural Network." *IEEE Computer* 21(3): 77-88.
- Cherkassky, V., and F. Mulier. 1998. *Learning from Data: Concepts, Theory, and Methods*. New York: Wiley.
- DeSieno, D. 1988. "Adding a Conscience Mechanism to Competitive Learning." In *IEEE International Conference on Neural Networks*, 117-124. Piscataway,

- NJ: IEEE Press.
- Feldman, J. A., and D. H. Ballard. 1982. "Connectionist Models and their Properties." *Cognitive Science* 6: 205-254.
- Fritzke, B. 1995. "Growing Cell Structures: A Self Organizing Network for Un-supervised and Supervised Training." *Neural Networks* 7: 1441-1460.
- Grossberg, S. 1980. "How does the Brain Build a Cognitive Code?" *Psychological Review* 87: 1-51.
- Hertz, J., A. Krogh, and R. G. Palmer. 1991. *Introduction to the Theory of Neural Computation*. Reading, MA: Addison Wesley.
- Jacobs, R. A., M. I. Jordan, S. J. Nowlan, and G. E. Hinton. 1991. "Adaptive Mixtures of Local Experts." *Neural Computation* 3: 79-87.
- Jordan, M. I., and R. A. Jacobs. 1994. "Hierarchical Mixtures of Experts and the EM Algorithm." *Neural Computation* 6: 181-214.
- Kohonen, T. 1990. "The Self-Organizing Map." *Proceedings of the IEEE* 78: 1464-1480.
- Kohonen, T. 1995. *Self-Organizing Maps*. Berlin: Springer.
- Lee, Y. 1991. "Handwritten Digit Recognition Using k -Nearest Neighbor, Radial Basis Function, and Backpropagation Neural Networks." *Neural Computation* 3: 440-449.
- Mao, J., and A. K. Jain. 1995. "Artificial Neural Networks for Feature Extraction and Multivariate Data Projection." *IEEE Transactions on Neural Networks* 6: 296-317.
- Moody, J., and C. Darken. 1989. "Fast Learning in Networks of Locally-Tuned Processing Units." *Neural Computation* 1: 281-294.
- Oja, E. 1982. "A Simplified Neuron Model as a Principal Component Analyzer." *Journal of Mathematical Biology* 15: 267-273.
- Omohundro, S. M. 1987. "Efficient Algorithms with Neural Network Behavior." *Complex Systems* 1: 273-347.
- Platt, J. 1991. "A Resource Allocating Network for Function Interpolation." *Neural Computation* 3: 213-225.
- Specht, D. F. 1991. "A General Regression Neural Network." *IEEE Transactions on Neural Networks* 2: 568-576.
- Tresp, V., J. Hollatz, and S. Ahmad. 1997. "Representing Probabilistic Rules with Networks of Gaussian Basis Functions." *Machine Learning* 27: 173-200.

核 机 器

核机器是最大边缘方法，允许把模型表示为训练实例的一个子集的影响之和。这些影响用应用专用的相似性核给出。我们讨论“核化的”分类、回归、离群点检测和维归约，以及如何选择和使用核。

13.1 引言

现在，我们讨论一种不同的线性分类和回归方法。不必惊奇，即使对于线性模型这种简单情况，也存在这么多不同的方法。每种学习算法都具有不同的归纳偏倚，做不同的假设，定义不同的目标函数，因此可能找到不同的线性模型。

本章，我们将讨论的模型称作支持向量机(Support Vector Machine, SVM)，稍后推广到核机器(kernel machine)。近年来它们非常流行，原因如下：

1) 它是基于判别式的方法，并使用 Vapnik 原则：不要在解决实际问题之前把解决一个更复杂的问题作为第一步(Vapnik 1995)。例如，对于分类，当任务是学习一个判别式时，不必估计类密度 $p(C_i | \mathbf{x})$ 或准确的后验概率值 $P(C_i | \mathbf{x})$ ；我们只需要估计类边界在哪里，即使 $P(C_i | \mathbf{x}) = P(C_j | \mathbf{x})$ 成立的 \mathbf{x} 。类似地，对于离群点检测，我们不需要估计全密度 $p(\mathbf{x})$ ；只需要找出把具有较低 $p(\mathbf{x})$ 值的 \mathbf{x} 分开的边界；即对于某个阈值 $\theta \in (0, 1)$ ，找出把满足 $p(\mathbf{x}) < \theta$ 的 \mathbf{x} 分开的边界。

2) 训练后，线性模型的参数(权重向量)可以用训练集的一个子集表示，这个子集因而称作支持向量(support vector)。对于分类，这些是靠近边界的实例，因此知道它们可以提取知识：这些是在两个类之间的边界附近、不确定或有错误的实例。它们的个数给我们提供了泛化误差的估计，并且正如我们将在下面看到的，用实例集表示模型参数可以进行核化(kernelization)。

3) 正如我们稍后将看到的，输出用支持向量的影响和表示，并且这些用核函数(kernel function)给出；核函数是数据实例之间相似性的应用专用度量。前面，我们谈到非线性基函数，它使得我们能够把输入映射到另一个空间，在那个空间里可以找到线性(光滑的)解。核函数使用相同的思想。

4) 通常，在大部分学习算法中，数据点用向量表示，或者使用点积(如在多层感知器中)，或者使用欧氏距离(如在径向基函数网络中)，核函数使得我们走得更远。例如， G_1 和 G_2 可以是两个图，而 $K(G_1, G_2)$ 可以对应于共享路径数，我们可以计算它，而不必显式地用向量表示 G_1 或 G_2 。

5) 基于核的算法可以形式化地表示成凸优化问题，并且存在单个最优解，我们可以解析地求解这个最优解。因此，我们不再受启发式方法的学习率、初始化、检查收敛性等困扰。当然，这并不意味着我们没有用于模型选择的超参数。我们有超参数，任何方法都需要它们，以使得算法与当前数据相匹配。

我们用分类开始我们的讨论, 然后推广到回归、离群点(新颖性)检测, 然后是维归约。我们看到, 在所有情况下, 基本上都有类似的二次规划问题, 在解的光滑性约束下, 最大化实例的分离性或边缘(margin)。通过对它求解我们得到支持向量。核函数根据其相似性概念定义空间, 并且一个核函数是好的, 如果在其对应空间中我们有更好的分离性。

309
310

13.2 最佳分离超平面

让我们还是从两类开始, 并使用 -1 和 $+1$ 标记这两个类。样本为 $\mathcal{X} = \{\mathbf{x}^t, r^t\}$, 其中如果 $\mathbf{x}^t \in C_1$ 则 $r^t = +1$, 如果 $\mathbf{x}^t \in C_2$ 则 $r^t = -1$ 。我们希望找到 \mathbf{w} 和 w_0 , 使得

$$\text{对于 } r^t = +1, \quad \mathbf{w}^T \mathbf{x}^t + w_0 \geq +1$$

$$\text{对于 } r^t = -1, \quad \mathbf{w}^T \mathbf{x}^t + w_0 \leq -1$$

它们可以改写为

$$r^t(\mathbf{w}^T \mathbf{x}^t + w_0) \geq +1 \quad (13.1)$$

注意, 我们并不是简单地要求

$$r^t(\mathbf{w}^T \mathbf{x}^t + w_0) \geq 0$$

为了更好地泛化, 我们不仅希望实例在超平面的正确一侧, 而且还希望它们离超平面有一定距离。超平面到它两侧最近实例的距离称作边缘(margin)。为了更好地泛化, 我们希望最大化边缘。

在 2.1 节讨论拟合一个矩形时谈到过边缘的概念, 我们指出最好把矩形放在 S 和 G 的中间, 留有余地, 这样做是为了在噪声少许移动检验实例时, 实例仍然在边界的正确一侧。

类似地, 现在我们使用的假设类是直线, 最佳分离超平面(optimal separating hyperplane)是最大化边缘的超平面。

回忆 10.3 节, \mathbf{x}^t 到判别式的距离为

$$\frac{|\mathbf{w}^T \mathbf{x}^t + w_0|}{\|\mathbf{w}\|}$$

当 $r^t \in \{-1, +1\}$ 时, 上式可以记作

$$\frac{r^t(\mathbf{w}^T \mathbf{x}^t + w_0)}{\|\mathbf{w}\|}$$

并且我们希望至少对于某个 ρ ,

$$\frac{r^t(\mathbf{w}^T \mathbf{x}^t + w_0)}{\|\mathbf{w}\|} \geq \rho, \quad \forall t \quad (13.2) \quad [311]$$

我们希望最大化 ρ , 但是缩放 \mathbf{w} , 我们可以得到的解有无限多个。为了得到唯一解, 我们固定 $\rho \|\mathbf{w}\| = 1$, 这样, 为最大化边缘, 我们最小化 $\|\mathbf{w}\|$, 因此这个任务可以定义为(参见 Cortes 和 Vapnik 1995; Vapnik 1995):

$$\min \frac{1}{2} \|\mathbf{w}\|^2, \quad \text{受限于 } r^t(\mathbf{w}^T \mathbf{x}^t + w_0) \geq +1, \quad \forall t \quad (13.3)$$

这是一个标准的二次优化问题, 其复杂度依赖于 d , 并且可以直接求解, 以找到 \mathbf{w} 和 w_0 。于是, 在超平面的两侧, 实例离超平面至少为 $1/\|\mathbf{w}\|$, 而整个边缘为 $2/\|\mathbf{w}\|$ 。

在 10.2 节中我们看到, 如果问题不是线性可分的, 则我们不拟合非线性函数, 而是使

用非线性基函数将问题映射到新的空间。通常,新空间的维度比原空间更高,并且在这种情况下,我们对复杂度不依赖输入维度的方法感兴趣。

在寻找最佳分离超平面时,我们可以把该优化问题转换成复杂度依赖于训练实例数 N ,而不依赖于 d 的形式。正如我们将在 13.5 节看到的,这种新表示方法的另一个优点是,它使得我们可以用核函数改写基函数。

为了得到新的公式,我们使用拉格朗日乘子 α^t ,将式(13.3)改写成非约束问题:

$$\begin{aligned} L_p &= \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{t=1}^N \alpha^t [r^t(\mathbf{w}^T \mathbf{x}^t + w_0) - 1] \\ &= \frac{1}{2} \|\mathbf{w}\|^2 - \sum_t \alpha^t r^t(\mathbf{w}^T \mathbf{x}^t + w_0) + \sum_t \alpha^t \end{aligned} \quad (13.4)$$

这应当关于 \mathbf{w} 、 w_0 最小化,并关于 $\alpha^t \geq 0$ 最大化。鞍点给出解。

这是一个凸二次优化问题,因为主要项是凸的,并且线性约束也是凸的。这样,我们可以使用 Karush-Kuhn-Tucker 条件,解其对偶问题。对偶问题是关于 α^t 最大化 L_p ,受限于约束 L_p 关于 \mathbf{w} 和 w_0 的梯度为 0,并且 $\alpha^t \geq 0$:

$$\frac{\partial L_p}{\partial \mathbf{w}} = 0 \Rightarrow \mathbf{w} = \sum_t \alpha^t r^t \mathbf{x}^t \quad (13.5)$$

$$\frac{\partial L_p}{\partial w_0} = 0 \Rightarrow \sum_t \alpha^t r^t = 0 \quad (13.6)$$

将它们代入式(13.4),我们得到对偶问题

$$\begin{aligned} L_d &= \frac{1}{2}(\mathbf{w}^T \mathbf{w}) - \mathbf{w}^T \sum_t \alpha^t r^t \mathbf{x}^t - w_0 \sum_t \alpha^t r^t + \sum_t \alpha^t \\ &= -\frac{1}{2}(\mathbf{w}^T \mathbf{w}) + \sum_t \alpha^t \\ &= -\frac{1}{2} \sum_t \sum_s \alpha^t \alpha^s r^t r^s (\mathbf{x}^t)^T \mathbf{x}^s + \sum_t \alpha^t \end{aligned} \quad (13.7)$$

我们只需要关于 α^t 对它最大化,受限于约束

$$\sum_t \alpha^t r^t = 0, \quad \text{并且对于任意 } t, \quad \alpha^t \geq 0$$

这可以使用二次优化方法来求解。对偶问题的规模依赖于样本的大小 N ,而不依赖于输入的维度 d 。时间复杂度的上界为 $O(N^3)$,而空间复杂度的上界为 $O(N^2)$ 。

一旦我们解出 α^t ,可以看到尽管它们有 N 个,但是多半以 $\alpha^t = 0$ 消失,而只有少量满足 $\alpha^t > 0$ 。 $\alpha^t > 0$ 的 \mathbf{x}^t 的集合是支持向量(support vector),并且正如我们在式(13.5)中看到的, \mathbf{w} 可以写成那些选作支持向量的训练实例的加权和。这些 \mathbf{x}^t 满足

$$r^t(\mathbf{w}^T \mathbf{x}^t + w_0) = 1$$

并且落在边缘上。我们使用这一事实,由任意支持向量来计算 w_0 :

$$w_0 = r^t - \mathbf{w}^T \mathbf{x}^t \quad (13.8)$$

从数值稳定性来讲,建议对所有支持向量计算上式,并取平均值。这样找出的判别式称作支持向量机(support vector machine, SVM)(参见图 13-1)。

α^t 大部分为 0,对于它们, $r^t(\mathbf{w}^T \mathbf{x}^t + w_0) > 1$ 。这些 \mathbf{x}^t 落在远离判别式的地方,并且对超平面没有影响。非支持向量的实例不携带信息,即便删除它们的任意子集,我们仍然得到

相同的解。从这个角度讲, SVM 算法与精简最近邻算法(参见 8.5 节)类似, 它只保存邻近(因而约束)类别别式的实例。

作为基于判别式的算法, SVM 只关注那些靠近边界的实例, 而丢弃那些落在内部的实例。使用这种思想, 可以在求 SVM 之前先使用一种较简单的分类器过滤掉这种实例的大部分, 从而降低 SVM 优化阶段的复杂度(习题 1)。

在检验阶段, 我们不强调边缘。我们计算 $g(\mathbf{x}) = \mathbf{w}^T \mathbf{x}' + w_0$, 并根据 $g(\mathbf{x})$ 的符号选择: 如果 $g(\mathbf{x}) > 0$ 则选择 C_1 , 否则选择 C_2 。

13.3 不可分情况: 软边缘超平面

如果数据不是线性可分的, 则我们前面讨论的算法就不能解决问题。在这种情况下, 如果两个类不是线性可分的, 使得不存在将它们分开的超平面, 则寻找出错最少的超平面。我们定义松弛变量(slack variable) $\xi^i \geq 0$, 存放到边缘的离差。有两种类型的离差: 一个实例可能位于超平面的错误一侧, 并被错误地分类; 或者实例可能在正确的一侧但落在边缘中, 即离超平面不够远。放宽式(13.1), 我们要求

$$r^i(\mathbf{w}^T \mathbf{x}' + w_0) \geq 1 - \xi^i \quad (13.9)$$

如果 $\xi^i = 0$, 则 \mathbf{x}' 没有问题。如果 $0 < \xi^i < 1$, 则 \mathbf{x}' 被正确分类, 但是它在边缘中。如果 $\xi^i \geq 1$, 则 \mathbf{x}' 被错误地分类(参见图 13-2)。错误分类数为 $\#\{\xi^i > 1\}$, 并且不可分的点数为 $\#\{\xi^i > 0\}$ 。我们定义软误差(soft error)为 $\sum_i \xi^i$, 并且将加上它作为罚项:

$$L_p = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i \xi^i \quad (13.10)$$

受限于式(13.9)的约束。 C 是罚因子, 像任意正则化模式一样, 在复杂度和数据误拟合之间权衡; 其中复杂度用权重向量的 L_2 范数度量(类似于多层感知器中的权衰减, 参见 11.9 节), 而数据误拟合用未分开的点数度量。注意, 为了更好地泛化, 我们不仅处罚误分类的点, 也处罚边缘中的点, 尽管后者在检验时被正确地分类。

加上这些约束, 式(13.4)的拉格朗日方程变成

$$L_p = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i \xi^i - \sum_i \alpha^i [r^i(\mathbf{w}^T \mathbf{x}' + w_0) - 1 + \xi^i] - \sum_i \mu^i \xi^i \quad (13.11)$$

其中 μ_i 是新的拉格朗日参数, 确保 ξ^i 为正。当我们对上式关于参数求导, 并令它们为 0 时, 我们得到

$$\frac{\partial L_p}{\partial \mathbf{w}} = \mathbf{w} - \sum_i \alpha^i r^i \mathbf{x}' = 0 \Rightarrow \mathbf{w} = \sum_i \alpha^i r^i \mathbf{x}' \quad (13.12)$$

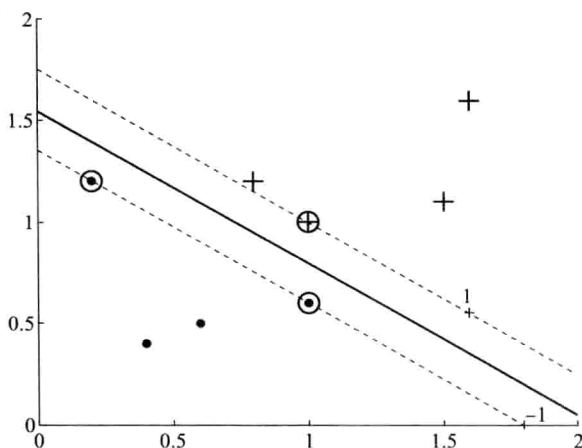


图 13-1 一个两类问题, 其中类实例用加号和圆点表示, 粗线是边界, 两侧的虚线定义边缘。圈住的实例是支持向量

$$\frac{\partial L_P}{\partial w_0} = \sum_t \alpha^t r^t = 0 \quad (13.13)$$

$$\frac{\partial L_P}{\partial \xi^t} = C - \alpha^t - \mu^t = 0 \quad (13.14)$$

由于 $\mu^t \geq 0$, 上式意味 $0 \leq \alpha^t \leq C$ 。把这些代入式(13.11), 我们得到关于 α^t 最大化的对偶问题

$$L_d = \sum_t \alpha^t - \frac{1}{2} \sum_t \sum_s \alpha^t \alpha^s r^t r^s (\mathbf{x}^t)^T \mathbf{x}^s \quad (13.15)$$

受限于

316

$$\sum_t \alpha^t r^t = 0, \quad \text{并且对于任意 } t, \quad 0 \leq \alpha^t \leq C$$

解这个对偶问题, 我们看到与可分情况一样, 落在边界正确一侧并距边界足够远的实例随 $\alpha^t = 0$ 消失(参见图 13-2)。支持向量的 $\alpha^t > 0$, 并且如式(13.12)所示, 它们定义 \mathbf{w} 。当然, 那些使 $\alpha^t < C$ 的实例在边缘上, 并且我们使用它们计算 w_0 ; 它们有 $\xi^t = 0$, 并且满足 $r^t (\mathbf{w}^T \mathbf{x}^t + w_0) = 1$ 。同样, 最好在这些 w_0 的估计上取平均值。在边缘中或误分类的那些实例的 $\alpha^t = C$ 。

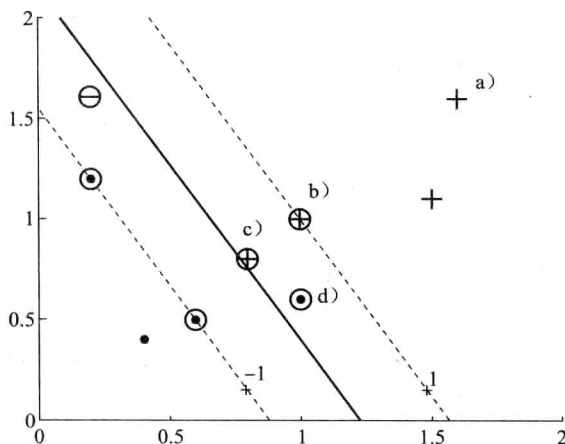


图 13-2 在对实例进行分类时, 有 4 种可能的情况: (a) 实例在正确一侧并且远离边缘; $r^t g(\mathbf{x}^t) > 1$, $\xi^t = 0$ 。(b) $\xi^t = 0$; 它在正确的一侧并且在边缘上。(c) $\xi^t = 1 - g(\mathbf{x})$, $0 < \xi^t \leq 1$; 点在正确的一侧, 但在边缘中, 离超平面不够远。(d) $\xi^t = 1 + g(\mathbf{x}) > 1$, 点在错误的一侧——这是误分类。除(a)之外, 所有实例都是支持向量。用对偶变量来说, 在(a)中有 $\alpha^t = 0$, 在(b)中有 $\alpha^t < C$, 在(c)和(d)中有 $\alpha^t = C$

作为支持向量存放的那些不可分的实例是这种实例, 即如果它们不在训练集中, 则我们会有麻烦, 不能正确地对它们分类; 它们要么被错误分类, 要么被正确分类但没有足够的置信度。我们可以说, 支持向量的个数是期望错误估计数的一个上界。并且实际上 Vapnik (1995) 已经证明期望检验错误率是

$$E_N[P(\text{error})] \leq \frac{E_N[\text{支持向量数}]}{N}$$

其中 $E_N[\cdot]$ 表示在大小为 N 的训练集上的期望。这意味着错误率依赖于支持向量的个数，而不依赖于输入的维度。

式(13.9)意味，如果实例在错误一侧，或者离边界的距离小于1，则我们定义它为错误，这称作转折点损失(hinge loss)。如果 $y^i = \mathbf{w}^T \mathbf{x}^i + w_0$ 是输出， r^i 是期望的输出，则转折点损失定义为

$$L_{\text{hinge}}(y^i, r^i) = \begin{cases} 0 & \text{如果 } y^i r^i \geq 1 \\ 1 - y^i r^i & \text{否则} \end{cases} \quad (13.16)$$

在图13-3中，我们把转折点损失与0/1损失、平方误差和互熵进行比较。我们看到，与0/1损失不同，转折点损失还处罚在边缘内的实例，尽管它们可能在正确的一侧，并且损失随着实例在错误一侧远离而线性增加。这也不同于平方损失，因而平方损失不如转折点损失鲁棒。我们看到互熵最小化逻辑斯谛判别式(参见10.7节)，或通过线性感知器(参见11.3节)，是对转折点损失的一个好的连续近似。

式(13.10)中的 C 是用交叉验证微调的正则化参数。它定义边缘最大化与误差最小化之间的权衡：如果它太大，则对未分开的点有高的处罚，并且可能存放许多支持向量而过拟合。如果它太小，则可能找到过于简单的解而欠拟合。通常，通过考察验证集上的准确率，在对数尺度从 $[10^{-6}, 10^{-5}, \dots, 10^{-5}, 10^{-6}]$ 中选择。

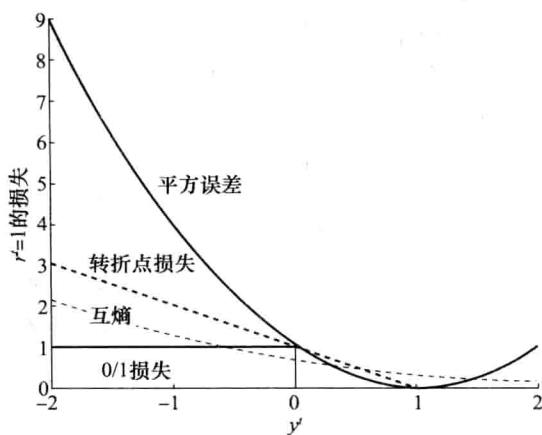


图13-3 比较 $r^i = 1$ 的不同损失函数：如果 $y^i = 1$ 则0/1损失为0，否则为1。如果 $y^i > 1$ 则转折点损失为0，否则为 $1 - y^i$ 。平方误差是 $(1 - y^i)^2$ 。互熵是 $\log(1/(1 + \exp(-y^i)))$

13.4 ν -SVM

另一种等价的软边缘超平面表示使用参数 $\nu \in [0, 1]$ ，而不是 C (Schölkopf 等, 2000)。目标函数是

$$\min \frac{1}{2} \|\mathbf{w}\|^2 - \nu \rho + \frac{1}{N} \sum_i \xi^i \quad (13.17)$$

受限于

$$r^i (\mathbf{w}^T \mathbf{x}^i + w_0) \geq \rho - \xi^i, \quad \xi^i \geq 0, \quad \rho \geq 0 \quad (13.18)$$

ρ 是一个新参数，它是优化问题的变量并且缩放边缘：现在边缘是 $2\rho/\|\mathbf{w}\|$ 。 ν 已经被证明是支持向量所占比例的一个下界和具有边缘误差 ($\sum_i \#\{\xi^i > 0\}$) 的实例所占比例的一个上界。对偶问题是

$$L_d = -\frac{1}{2} \sum_i \sum_s \alpha^i \alpha^s r^i r^s (\mathbf{x}^i)^T \mathbf{x}^s \quad (13.19)$$

受限于

$$\sum_i \alpha' r^i = 0, \quad 0 \leq \alpha' \leq \frac{1}{N}, \quad \sum_i \alpha' \geq v$$

当我们将式(13.19)与式(13.15)比较时, 可以看到项 $\sum_i \alpha'$ 不再出现在目标函数中, 而是一个常数。通过操控 v , 我们可以控制支持向量所占的比例, 并且这被认为比操控 C 更直观。

13.5 核技巧

10.2 节指出, 如果问题是非线性的, 则我们不是拟合一个非线性模型, 而是使用合适的基函数, 通过非线性变换将问题映射到新空间, 然后在新空间使用线性模型。新空间中的线性模型对应于原空间中的非线性模型。这种方法可以用于分类和回归问题, 并且对于分类这种特殊情况, 它可以与任意模式共同使用。在支持向量机这种特定情况下, 这导致某种简化。现在我们来讨论。

设我们有通过基函数

$$\mathbf{z} = \boldsymbol{\phi}(\mathbf{x}), \quad \text{其中 } z_j = \phi_j(\mathbf{x}), \quad j = 1, \dots, k$$

计算新的维度, 将 d -维 \mathbf{x} 空间映射到 k -维 \mathbf{z} 空间, 其中判别式表示为

$$\begin{aligned} g(\mathbf{z}) &= \mathbf{w}^T \mathbf{z} \\ g(\mathbf{x}) &= \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}) \\ &= \sum_{j=1}^k w_j \phi_j(\mathbf{x}) \end{aligned} \quad (13.20)$$

这里, 我们并不使用单独的 w_0 ; 我们假定 $z_1 = \phi_1(\mathbf{x}) \equiv 1$ 。通常, k 远大于 d , k 也可能大于 N , 并且这也是使用对偶形式的优点, 对偶问题的复杂度依赖于 N , 如果我们使用原问题, 则将依赖于 k 。这里, 我们也使用软边缘超平面更一般的情况, 因为我们不能保证问题在新空间内是线性可分的。

319

问题是一样的

$$L_p = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i \xi^i \quad (13.21)$$

不同之处是现在约束定义在新空间

$$r^i \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}^i) \geq 1 - \xi^i \quad (13.22)$$

拉格朗日方程是

$$L_p = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i \xi^i - \sum_i \alpha^i [r^i \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}^i) - 1 + \xi^i] - \sum_i \mu^i \xi^i \quad (13.23)$$

当我们关于参数求导并令它们等于 0 时, 得到

$$\frac{\partial L_p}{\partial \mathbf{w}} = \mathbf{w} = \sum_i \alpha^i r^i \boldsymbol{\phi}(\mathbf{x}^i) \quad (13.24)$$

$$\frac{\partial L_p}{\partial \xi^i} = C - \alpha^i - \mu^i = 0 \quad (13.25)$$

现在, 对偶问题是

$$L_d = \sum_i \alpha^i - \frac{1}{2} \sum_i \sum_s \alpha^i \alpha^s r^i r^s \boldsymbol{\phi}(\mathbf{x}^i)^T \boldsymbol{\phi}(\mathbf{x}^s) \quad (13.26)$$

受限于

$$\sum_t \alpha^t r^t = 0 \text{ 并且对于任意 } t, \quad 0 \leq \alpha^t \leq C$$

核机器的基本思想是用支持向量和原输入空间中的实例之间的核函数 $K(\mathbf{x}', \mathbf{x}')$ 取代基函数的内积 $\phi(\mathbf{x}')^T \phi(\mathbf{x}')$ 。这样, 不是把两个实例 \mathbf{x}' 和 \mathbf{x}^s 映射到 z 空间并在那里做点积, 而是直接使用原空间中的核函数。

$$L_d = \sum_t \alpha^t - \frac{1}{2} \sum_t \sum_s \alpha^t \alpha^s r^t r^s K(\mathbf{x}', \mathbf{x}^s) \quad (13.27) \quad \boxed{320}$$

核函数也出现在判别式中

$$\begin{aligned} g(\mathbf{x}) &= \mathbf{w}^T \phi(\mathbf{x}) = \sum_t \alpha^t r^t \phi(\mathbf{x}') \phi(\mathbf{x}) \\ &= \sum_t \alpha^t r^t K(\mathbf{x}', \mathbf{x}) \end{aligned} \quad (13.28)$$

这意味着, 如果我们有核函数, 则完全不需要把它映射到新空间。实际上, 对于任何有效的核函数, 确实存在对应的映射函数, 但是使用 $K(\mathbf{x}', \mathbf{x})$ 比计算 $\phi(\mathbf{x}')$ 和 $\phi(\mathbf{x})$, 再求点积简单得多。正如我们在下一节所看到的, 许多算法都被核化, 而这正是我们称为“核机器”的原因。

核值的矩阵 \mathbf{K} , 其中 $\mathbf{K}_{ts} = K(\mathbf{x}', \mathbf{x}^s)$, 称作 Gram 矩阵 (Gram matrix), 该矩阵应该是对称的、半正定的。近来, 在共享数据集中, 仅有 \mathbf{K} 矩阵而不提供 \mathbf{x}' 或 $\phi(\mathbf{x}')$ 已经成为标准做法。尤其是在生物信息学和自然语言处理的应用中, \mathbf{x} (或 $\phi(\mathbf{x})$) 有数百维或数千维, 存放或下载这个 $N \times N$ 矩阵开销小得多 (Vert、Tsuda 和 Schölkopf 2004); 然而, 这意味着我们只能使用这些可用的信息进行训练或检验, 并且不能使用训练后的模型对该数据集之外的数据进行预测。

13.6 向量核

最流行的通用核函数是

- q 次多项式:

$$K(\mathbf{x}', \mathbf{x}) = (\mathbf{x}^T \mathbf{x}^t + 1)^q \quad (13.29)$$

其中 q 由用户选择。例如, 当 $q = 2$, $d = 2$ 时,

$$\begin{aligned} K(\mathbf{x}, \mathbf{y}) &= (\mathbf{x}^T \mathbf{y} + 1)^2 \\ &= (x_1 y_1 + x_2 y_2 + 1)^2 \\ &= 1 + 2x_1 y_1 + 2x_2 y_2 + 2x_1 x_2 y_1 y_2 \\ &\quad + x_1^2 y_1^2 + x_2^2 y_2^2 \end{aligned}$$

它对应于如下基函数的内积 (Cherkassky 和 Mulier 1998):

$$\phi(\mathbf{x}) = [1, \sqrt{2}x_1, \sqrt{2}x_2, \sqrt{2}x_1 x_2, x_1^2, x_2^2]^T$$

图 13-4 中给出了一个例子。当 $q = 1$ 时, 我们有对应于原公式的线性核 (linear kernel)。

- 径向基函数 (radial-based function):

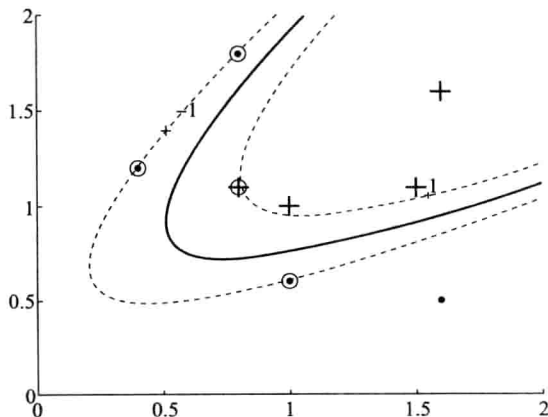


图 13-4 被二次多项式核找出的决策边界和边缘。
圈住的实例是支持向量

$$K(\mathbf{x}', \mathbf{x}) = \exp \left[-\frac{\|\mathbf{x}' - \mathbf{x}\|^2}{2s^2} \right] \quad (13.30)$$

与 Parzen 窗口一样(参见第 8 章), 它定义一个球形核, 其中 \mathbf{x}' 是中心, 而 s 由用户提供, 定义半径, 这也类似于第 12 章讨论的径向基函数。

图 13-5 中给出了一个例子。从这个例子中我们看到, 较大的展览光滑了边界; 最好的值用交叉验证找到。注意, 当有两个参数使用交叉验证优化时(例如这里的 C 和 s^2), 应该在两维上进行栅格(因子)搜索。我们将在 19.2 节讨论搜索这种因子最佳组合的方法。

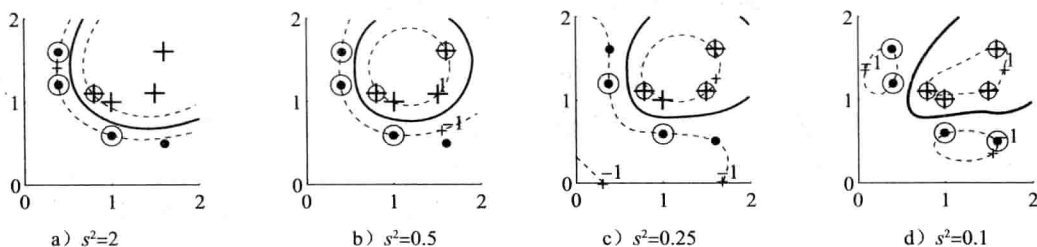


图 13-5 被具有不同展览值 s^2 的高斯核找出的决策边界和边缘。使用较大展览, 我们得到更光滑的边界

我们可以推广欧氏距离, 从而使用马氏距离核:

$$K(\mathbf{x}', \mathbf{x}) = \exp \left[-\frac{1}{2}(\mathbf{x}' - \mathbf{x})^T \mathbf{S}^{-1}(\mathbf{x}' - \mathbf{x}) \right] \quad (13.31)$$

其中 \mathbf{S} 是协方差矩阵。或者, 在最一般的情况下, 对于某个距离函数 $\mathcal{D}(\mathbf{x}', \mathbf{x})$, 使用

$$K(\mathbf{x}', \mathbf{x}) = \exp \left[-\frac{\mathcal{D}(\mathbf{x}', \mathbf{x})}{2s^2} \right] \quad (13.32)$$

■ S 形函数(sigmoidal function):

$$K(\mathbf{x}', \mathbf{x}) = \tanh(2\mathbf{x}^T \mathbf{x}' + 1) \quad (13.33)$$

其中 $\tanh(\cdot)$ 与 S 形函数具有相同的形状, 不同的是它的取值为 $-1 \sim +1$, 这类似于我们在第 11 章讨论的多层感知器。

13.7 定义核

还可以定义应用专用的核。核通常被看作相似性度量, 意指从应用的角度来看, 当 \mathbf{x} 和 \mathbf{y} 更“相似”时, $K(\mathbf{x}, \mathbf{y})$ 取更大的值。这意味着, 关于应用的任何先验知识都可以通过定义合适的核提供给学习算法(“核工程”), 并且核的这种用法可以看作“线索”(参见 11.8.4 节)的另一个例子。

有串核、树核、图核等(Vert、Tsuda 和 Schölkopf 2004), 取决于我们如何表示数据, 如何度量这种表示下的相似性。

例如, 给定两份文档, 出现在这两份文档中词的个数可以作为核。假定 D_1 和 D_2 是两份文档, 而一种可能的表示称作词袋(bag of words), 其中我们预先定义了 M 个与应用相关的

词。定义 $\phi(D_1)$ 为 M 维二元向量, 如果第 i 个词出现在 D_1 中, 则其第 i 个分量为 1, 否则为 0。于是, $\phi(D_1)^T \phi(D_2)$ 计算被共享的词个数。这里, 我们看到, 如果直接把 $K(D_1, D_2)$ 作为共享词的个数来定义和实现, 则我们不需要预先选择 M 个词, 而是使用词典中的任何词(当然, 在丢弃诸如 of、and 等不提供信息的词之后), 并且不需要明确地产生词袋表示, 仿佛我们允许 M 想多大就多大一样。

有时(例如在生物信息学应用中), 我们可以计算两个对象之间的相似度得分(similarity score), 这不必是半正定的。给定两个(基因)串, 一种核度量是编辑距离(edit distance), 即把一个串转换成另一个需要做多少次操作(插入、删除、替换), 这又称为比对(alignment)。在这种情况下, 一种技巧是定义一个 M 个模板的集合, 并且把对象表示成一个在所有模板上得分的 M 维向量。也就是说, 如果 $m_i (i=1, \dots, M)$ 是模板, $s(\mathbf{x}', m_i)$ 是 \mathbf{x}' 和 m_i 之间的得分, 则我们定义

$$\phi(\mathbf{x}') = [s(\mathbf{x}', m_1), s(\mathbf{x}', m_2), \dots, s(\mathbf{x}', m_M)]^T$$

并且我们定义经验核映射(empirical kernel map)为

$$K(\mathbf{x}', \mathbf{x}^s) = \phi(\mathbf{x}')^T \phi(\mathbf{x}^s)$$

323
?
324

这是一个合法的核。

有时, 我们有二元评分函数, 例如, 两种蛋白质可能相互作用, 也可能不。我们希望把这一点推广到任意两个二元实例的得分。在这种情况下, 技巧是定义一个图, 其中节点是实例, 如果它们相互作用, 即如果该二元评分返回 1, 而两个节点被连接。于是, 如果它们之间的路径短或被许多路径连接, 则认为两个不直接连接的节点是“相似的”。这把逐对局部相互作用转换成全局相似性度量, 很像 Isomap 使用的测地距离(参见 6.7 节), 并称为扩散核(diffusion kernel)。

如果 $p(\mathbf{x})$ 是概率密度, 则

$$K(\mathbf{x}', \mathbf{x}) = p(\mathbf{x}')p(\mathbf{x})$$

是一个合法的核函数。当 $p(\mathbf{x})$ 是 \mathbf{x} 的生成模型, 用于度量我们看到 \mathbf{x} 的可能性时, 使用这个核函数。例如, 如果 \mathbf{x} 是一个序列, 则 $p(\mathbf{x})$ 可能是隐马尔可夫模型(参见第 15 章)。使用这个核函数, 如果 \mathbf{x}' 和 \mathbf{x} 可能都是被相同模型生成的, 则 $K(\mathbf{x}', \mathbf{x})$ 将取很高的值。还可以将生成模型参数化为 $p(\mathbf{x} | \theta)$, 并且由数据学习 θ , 这称作费希尔核(Fisher kernel)(Jaakkola 和 Haussler 1998)。

13.8 多核学习

可以通过组合一些简单的核来构造新的核。如果 $K_1(\mathbf{x}, \mathbf{y})$ 和 $K_2(\mathbf{x}, \mathbf{y})$ 是两个合法的核, 而 c 是常数, 则

$$K(\mathbf{x}, \mathbf{y}) = \begin{cases} cK_1(\mathbf{x}, \mathbf{y}) \\ K_1(\mathbf{x}, \mathbf{y}) + K_2(\mathbf{x}, \mathbf{y}) \\ K_1(\mathbf{x}, \mathbf{y}) \cdot K_2(\mathbf{x}, \mathbf{y}) \end{cases} \quad (13.34)$$

也是合法的核。

还可以对 \mathbf{x} 的不同子集使用不同的核。这样, 我们看到组合核可以作为融合来自不同信息源信息的另一种方法, 其中每个核都根据自己的领域度量相似性。当我们有来自两种表示

A 和 B 的输入时,

$$\begin{aligned} K_A(\mathbf{x}_A, \mathbf{y}_A) + K_B(\mathbf{x}_B, \mathbf{y}_B) &= \boldsymbol{\phi}_A(\mathbf{x}_A)^T \boldsymbol{\phi}_A(\mathbf{y}_A) + \boldsymbol{\phi}_B(\mathbf{x}_B)^T \boldsymbol{\phi}_B(\mathbf{y}_B) \\ &= \boldsymbol{\phi}(\mathbf{x})^T \boldsymbol{\phi}(\mathbf{y}) \\ &= K(\mathbf{x}, \mathbf{y}) \end{aligned} \quad (13.35)$$

其中 $\mathbf{x} = [\mathbf{x}_A, \mathbf{x}_B]$ 是两个表示的串接。也就是说, 取两个核的和对应于在串接特征向量中做点积, 这可以推广到大量核

$$K(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^m K_i(\mathbf{x}, \mathbf{y}) \quad (13.36)$$

这类似于取分类器的平均(参见 17.4 节), 这次是在核上取平均, 不再需要我们去选择一个特定的核。还可以取加权和, 并由数据学习权重(Lanckriet 等 2004, Sonnenburg 等 2006):

$$K(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^m \eta_i K_i(\mathbf{x}, \mathbf{y}) \quad (13.37)$$

满足 $\eta_i \geq 0$, 满足或不满足约束 $\sum_i \eta_i = 1$, 分别称作凸组合和锥形组合, 这称作多核学习(multiple kernel learning), 其中我们用核的加权和取代单核。单核目标函数式(13.27)变成

$$L_d = \sum_t \alpha^t - \frac{1}{2} \sum_t \sum_s \alpha^t \alpha^s r^t r^s \sum_i \eta_i K_i(\mathbf{x}^t, \mathbf{x}^s) \quad (13.38)$$

我们求解支持向量机参数 α^t 和核权重 η_i 。多核的组合也出现在判别式中

$$g(\mathbf{x}) = \sum_t \alpha^t r^t \sum_i \eta_i K_i(\mathbf{x}^t, \mathbf{x}) \quad (13.39)$$

训练后, η_i 的取值取决于对应的核 $K_i(\mathbf{x}^t, \mathbf{x})$ 在判别式中的作用。还可以通过把核权重定义为输入 \mathbf{x} 的有参函数对核的局部化, 这非常像混合专家模型的选通函数(参见 17.8 节)

$$g(\mathbf{x}) = \sum_t \alpha^t r^t \sum_i \eta_i(\mathbf{x} | \theta_i) K_i(\mathbf{x}^t, \mathbf{x}) \quad (13.40)$$

并且选通参数 θ_i 与支持向量机参数一起学习(Gönen 和 Alpaydin 2008)。

当我们有来自多个不同的表示或不同形态下的数据源信息时(例如, 在语音识别中我们可能有声波和唇动图像), 通常的方法是把它们分别提供给不同的分类器, 然后融合它们的决策; 我们将在第 17 章详细讨论这种方法。组合多个核提供了另一种集成多源输入的方法, 其中单个分类器对非同源的输入使用不同的核, 因此有不同的相似性概念(Noble 2004)。于是, 局部化版本可以看作它的扩展, 其中我们可以根据输入来选择数据源, 从而选择相似性度量。

13.9 多类核机器

当存在 $K > 2$ 个类时, 直截了当地, 一对所有(one-vs-all)方法定义 K 个两类问题, 每个都把一个类与其他所有类的组合分开, 并学习 K 个支持向量机 $g_i(\mathbf{x}) (i = 1, \dots, K)$ 。也就是说, 在训练 $g_i(\mathbf{x})$ 时, C_i 类的实例标记为 +1, 而 $C_k (k \neq i)$ 类的实例标记为 -1。在检验时, 我们计算所有的 $g_i(\mathbf{x})$, 并选择最大的。

Platt(1999)提出用一个 S 形函数拟合单个(两类)SVM 的输出, 以把输出转换成后验概率。类似地, 可以训练一个软最大输出层来最小化互熵, 以产生 $K > 2$ 个后验概率(Mayoraz 和 Alpaydin 1999):

$$y_i(\mathbf{x}) = \sum_{j=1}^K v_{ij} f_j(\mathbf{x}) + v_{i0} \quad (13.41)$$

其中 $f_j(\mathbf{x})$ 是 SVM 的输出, 而 y_i 是后验概率输出。训练权重 v_{ij} 来最小化互熵。然而需要注意, 与层叠中一样(参见 17.9 节), 训练 v_{ij} 的数据应当不同于用来训练基本支持向量机 $f_j(\mathbf{x})$ 的数据, 以便减轻过拟合。

与其他分类方法一样, 不像通常那样构建 K 个把一个类与其余类分开的两类 SVM 分类器, 而是构建 $K(K-1)/2$ 个逐对(pairwise)分类器(参见 10.4 节), 每个 $g_{ij}(\mathbf{x})$ 取 C_i 类实例标记为 +1, C_j 类实例标记为 -1, 并且不使用其他类的实例。一般认为逐对把类分开是一个较简单的事, 另外的优点是使用较少的数据, 优化更快; 但要注意, 这需要训练 $\mathcal{O}(K^2)$ 个而不是 $\mathcal{O}(K)$ 个判别式。

在一般情况下, 一对所有和逐对分开都是把一个多类问题分解成一组两类问题的校正输出码(error-correcting output codes)的特例(Dietterich 和 Bakiri 1995)(参见 17.6 节)。作为两类分类器的 SVM 是两类分类的理想选择(Allwein, Schapire 和 Singer 2000), 并且还可以有增量方法, 增加新的两类 SVM, 以更好地把存在问题的一对类分开, 从而改善不理想的 ECOC 矩阵(Mayoraz 和 Alpaydin 1998)。

另一种可能的方法是设计一个涉及所有类的多类(multiclass)优化问题(Weston 和 Watkins 1998):

$$\min \frac{1}{2} \sum_{i=1}^K \|\mathbf{w}_i\|^2 + C \sum_i \sum_t \xi_i^t \quad (13.42)$$

受限于

$$\text{对于任意的 } i \neq z^t, \quad \mathbf{w}_i \cdot \mathbf{x}^t + w_{z0} \geq \mathbf{w}_i \cdot \mathbf{x}^t + w_{z0} + 2 - \xi_i^t, \quad \text{并且 } \xi_i^t \geq 0$$

其中, z^t 包含 \mathbf{x}^t 的类索引。正则化项同时减少所有超平面的范数, 而约束则确保该类与任何其他类之间的边缘至少为 2。正确类的输出应当至少为 +1, 其他类的输出应当至少为 -1, 而松弛变量用来补差。

尽管这看上去漂亮, 但是一对所有方法通常更可取, 因为它分别解决 K 个 N 个变量的问题, 而多类方法使用 $K \cdot N$ 个变量。

13.10 用于回归的核机器

现在, 让我们看看如何将支持向量机推广到回归问题。我们看到相同的定义可接受边缘的方法、松弛变量, 综合光滑性和误差的正则化函数在这里也能用。我们从线性模型开始, 稍后我们关注如何使用核函数:

$$f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$$

对于一般的回归, 我们使用差的平方作为误差:

$$e_2(r^t, f(\mathbf{x}^t)) = [r^t - f(\mathbf{x}^t)]^2$$

然而, 对于支持向量回归, 我们使用 ε -敏感损失函数:

$$e_\varepsilon(r^t, f(\mathbf{x}^t)) = \begin{cases} 0 & \text{如果 } |r^t - f(\mathbf{x}^t)| < \varepsilon \\ |r^t - f(\mathbf{x}^t)| - \varepsilon & \text{否则} \end{cases} \quad (13.43)$$

这意味我们容忍高达 ε 的误差, 并且超出的误差具有线性而不是平方影响。因此, 这种误差

函数更能抵御噪声,因而更加鲁棒(参见图13-6)。与转折点损失一样,有一个区域没有误差,这导致了稀疏性。

类似于软边缘超平面,我们引入松弛变量来处理 ε -区域之外的偏差,并且得到(Vapnik 1995):

$$\min \frac{1}{2} \| \mathbf{w} \|^2 + C \sum_t (\xi_+^t + \xi_-^t) \quad (13.44)$$

受限于

$$r^t - (\mathbf{w}^T \mathbf{x} + w_0) \leq \varepsilon + \xi_+^t$$

$$(\mathbf{w}^T \mathbf{x} + w_0) - r^t \leq \varepsilon + \xi_-^t$$

$$\xi_+^t, \xi_-^t \geq 0$$

这里,对正和负的偏差,我们使用了两种类型的松弛变量,以保持它们为正。实际上,我们可以把这看作两个背靠背加上的转折点损失,一个用于正松弛,一个用于负松弛。该式对应于式(13.43)中给出的 ε -敏感损失函数。拉格朗日方程是

$$\begin{aligned} L_p &= \frac{1}{2} \| \mathbf{w} \|^2 + C \sum_t (\xi_+^t + \xi_-^t) \\ &\quad - \sum_t \alpha_+^t [\varepsilon + \xi_+^t - r^t + (\mathbf{w}^T \mathbf{x} + w_0)] \\ &\quad - \sum_t \alpha_-^t [\varepsilon + \xi_-^t - (\mathbf{w}^T \mathbf{x} + w_0) + r^t] - \sum_t (\mu_+^t \xi_+^t + \mu_-^t \xi_-^t) \end{aligned} \quad (13.45)$$

取偏导数,我们得到

$$\frac{\partial L_p}{\partial \mathbf{w}} = \mathbf{w} - \sum_t (\alpha_+^t - \alpha_-^t) \mathbf{x}^t = 0 \Rightarrow \mathbf{w} = \sum_t (\alpha_+^t - \alpha_-^t) \mathbf{x}^t \quad (13.46)$$

$$\frac{\partial L_p}{\partial w_0} = \sum_t (\alpha_+^t - \alpha_-^t) = 0 \quad (13.47)$$

$$\frac{\partial L_p}{\partial \xi_+^t} = C - \alpha_+^t - \mu_+^t = 0 \quad (13.48)$$

$$\frac{\partial L_p}{\partial \xi_-^t} = C - \alpha_-^t - \mu_-^t = 0 \quad (13.49)$$

对偶问题是

$$\begin{aligned} L_d &= -\frac{1}{2} \sum_t \sum_s (\alpha_+^t - \alpha_-^t) (\alpha_+^s - \alpha_-^s) (\mathbf{x}^t)^T \mathbf{x}^s \\ &\quad - \varepsilon \sum_t (\alpha_+^t + \alpha_-^t) + \sum_t r^t (\alpha_+^t - \alpha_-^t) \end{aligned} \quad (13.50)$$

受限于

$$0 \leq \alpha_+^t \leq C, \quad 0 \leq \alpha_-^t \leq C, \quad \sum_t (\alpha_+^t - \alpha_-^t) = 0$$

一旦我们求出它的解,就会看到落入管(tube)中的所有实例都有 $\alpha_+^t = \alpha_-^t = 0$; 这些是以足够精度拟合的实例(参见图13-7)。支持向量满足 $\alpha_+^t > 0$ 或 $\alpha_-^t > 0$, 并且都是这两种类型。它们可能是管边界上的实例(α_+^t 或 α_-^t 在 $0 \sim C$ 之间), 并且我们使用它们计算 w_0 。例

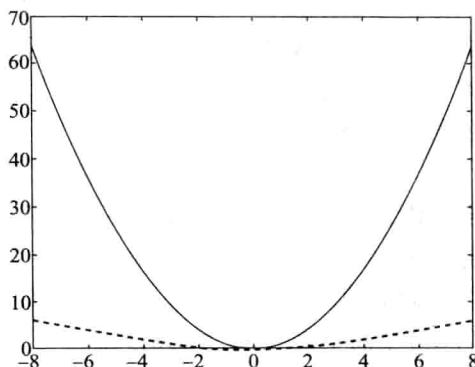


图13-6 平方误差函数和 ε -敏感误差函数。我们看到 ε -敏感误差函数不受小误差的影响,并且受大误差的影响也较小,因此对离群点更鲁棒

如, 假定 $\alpha'_+ > 0$, 我们有 $r' = \mathbf{x}'^T \mathbf{x}' + w_0 + \varepsilon$ 。落在 ε -管外的实例是第二种类型; 这些是没有很好拟合的实例 ($\alpha'_+ = C$), 如图 13-7 所示。

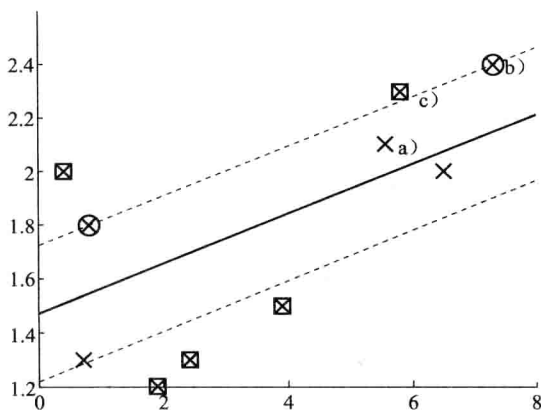


图 13-7 “x”表示数据点的拟合回归线, 并显示了 ε -管 ($C=10$, $\varepsilon=0.25$)。有 3 种情况: (a) 实例在管中; (b) 实例在管的边界上 (圆中的实例); (c) 实例在管外, 具有正的松弛, 即 $\xi'_+ > 0$ (方框中的实例)。(b) 和 (c) 是支持向量。用对偶变量来说, 在 (a) 中有 $\alpha'_+ = 0$, $\alpha'_- = 0$; 在 (b) 中有 $\alpha'_+ < C$; 在 (c) 中有 $\alpha'_+ = C$

使用式 (13.46), 我们可以把拟合直线写成支持向量的加权和:

$$f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0 = \sum_i (\alpha'_+ - \alpha'_-) (\mathbf{x}'^i)^T \mathbf{x} + w_0 \quad (13.51)$$

式 (13.50) 中的点积 $(\mathbf{x}'^i)^T \mathbf{x}$ 也可以用核函数 $K(\mathbf{x}'^i, \mathbf{x})$ 替换, 并且类似地, $(\mathbf{x}'^i)^T \mathbf{x}$ 可以用核函数 $K(\mathbf{x}'^i, \mathbf{x})$ 替换, 并得到非线性拟合。使用多项式核类似于拟合一个多项式 (参见图 13-8), 而使用高斯核 (参见图 13-9) 则类似于非参数光滑模型 (参见 8.6 节), 不同之处在于由于解的稀疏性, 我们不需要整个训练集, 而只需要一个子集。

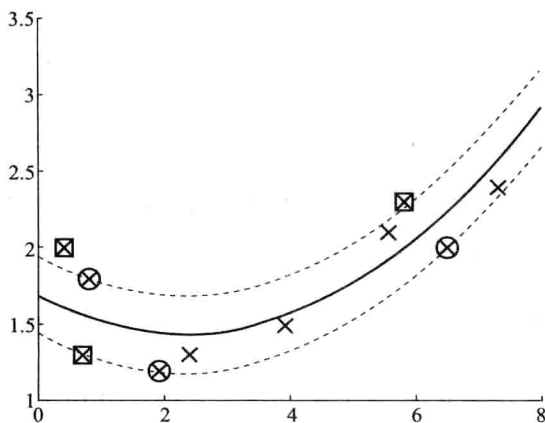


图 13-8 显示使用二次核函数拟合的回归线和 ε -管 ($C=10$, $\varepsilon=0.25$)。圆中的实例是边缘上的支持向量, 方框中的实例是离群点支持向量

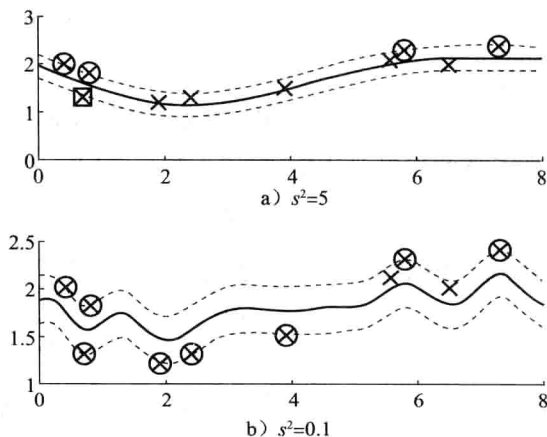


图 13-9 显示使用两个具有不同展宽的高斯核拟合的回归线和 ε -管 ($C=10$, $\varepsilon=0.25$)。圆中的实例是边缘上的支持向量, 方框中的实例是离群点支持向量

331
332

对于回归, 也有一个等价的 ν -SVM (Schölkopf 等 2000), 其中取代固定的 ε , 我们固定 ν 来限制支持向量的比例, 但仍然需要 C 。

13.11 一类核机器

支持向量机最初是为分类提出的, 通过为回归线而不是为判别式附近的偏差定义松弛变量, SVM 被扩展到回归。现在, 我们看看如何把 SVM 用于一类受限的非监督学习, 即估计高密度区域。我们并不进行整个密度估计, 而是想找出把高密度区域与低密度区域分开的边界 (因此它像一个分类问题) (Tax 和 Duin 1999)。这种边界可以用于新颖性 (novelty) 或离群点检测 (outlier detection), 这也称作一类分类 (one-class classification)。

考虑一个中心为 \mathbf{a} 、半径为 R 的球。我们希望它围住的密度尽可能大, 其中密度根据经验用训练集百分比度量。同时, 与之权衡, 我们希望找到最小半径 (参见图 13-10)。我们为落在球外的实例定义一个松弛变量 (只有一种类型的松弛变量, 因为我们只有来自一个类的实例, 并且对落在球内的那些实例没有处罚), 并且有一个正比于半径的光滑性度量:

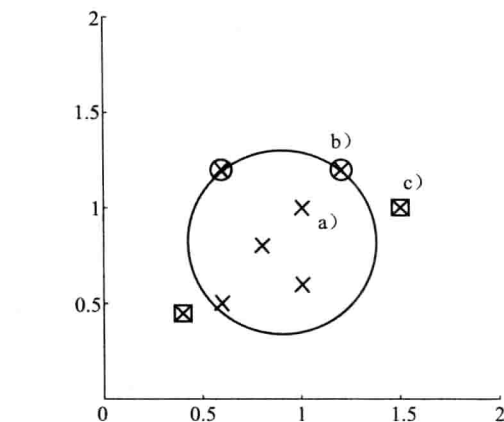


图 13-10 一类支持向量机把光滑的边界 (这里使用线性核函数, 圆具有最小半径) 放置在尽可能多地围住实例的地方。存在 3 种可能的情况: (a) 实例是典型实例; (b) $\xi^t = 0$ 的实例落在边界上, 这种实例定义 R ; (c) $\xi^t > 0$ 的实例是离群点。(b) 和 (c) 是支持向量。用对偶变量的术语来说, 在 (a) 中有 $\alpha^t = 0$; 在 (b) 中有 $0 < \alpha^t < C$; 在 (c) 中有 $\alpha^t = C$

$$\min R^2 + C \sum_t \xi^t \quad (13.52)$$

受限于

$$\|\mathbf{x}^t - \mathbf{a}\|^2 \leq R^2 + \xi^t, \quad \text{并且对于任意的 } t \text{ 有 } \xi^t \geq 0$$

加上这些约束, 注意到 $\|\mathbf{x}^t - \mathbf{a}\|^2 = (\mathbf{x}^t - \mathbf{a})^T (\mathbf{x}^t - \mathbf{a})$, 我们得到拉格朗日方程:

$$L_p = R^2 + C \sum_t \xi^t - \sum_t \alpha^t (R^2 + \xi^t - [(\mathbf{x}^t)^T \mathbf{x}^t - 2\mathbf{a}^T \mathbf{x}^t + \mathbf{a}^T \mathbf{a}]) - \sum_t \gamma^t \xi^t \quad (13.53)$$

其中 $\alpha^t \geq 0$, $\gamma^t \geq 0$ 是拉格朗日乘子。关于参数求导, 我们得到

$$\frac{\partial L_p}{\partial R} = 2R - 2R \sum_t \alpha^t = 0 \Rightarrow \sum_t \alpha^t = 1 \quad (13.54)$$

$$\frac{\partial L_p}{\partial \mathbf{a}} = \sum_t \alpha^t (2\mathbf{x}^t - 2\mathbf{a}) = 0 \Rightarrow \mathbf{a} = \sum_t \alpha^t \mathbf{x}^t \quad (13.55)$$

$$\frac{\partial L_p}{\partial \xi^t} = C - \alpha^t - \gamma^t = 0 \quad (13.56)$$

由于 $\gamma^t \geq 0$, 因此我们可以把最后一个约束写成: $0 \leq \alpha^t \leq C$ 。把这些代入式 (13.53), 我们得到关于 α^t 最大化的对偶问题:

333

$$L_d = \sum_i \alpha^i (\mathbf{x}^i)^T \mathbf{x}^i - \sum_i \sum_s \alpha^i \alpha^s (\mathbf{x}^i)^T \mathbf{x}^s \quad (13.57)$$

受限于

$$0 \leq \alpha^i \leq C \text{ 并且 } \sum_i \alpha^i = 1$$

当我们求解该优化问题时, 再次看到大部分实例随着它们的 $\alpha^i = 0$ 消失; 这些是落在球体内的典型、高度相似的实例(参见图 13-10)。有两种类型的支持向量满足 $\alpha^i > 0$: 满足 $0 < \alpha^i < C$ 并且落在边界上的实例, $\|\mathbf{x}^i - \mathbf{a}\|^2 = R^2$ ($\xi^i = 0$), 这些用来计算 R 。满足 $\alpha^i = C$ ($\xi^i > 0$) 的实例落在边界外, 是离群点。由式(13.55), 我们看到中心 \mathbf{a} 是支持向量的加权和。

于是, 给定一个检验输入 \mathbf{x} , 我们说它是离群点, 如果

$$\|\mathbf{x}^i - \mathbf{a}\|^2 > R^2$$

或

$$\mathbf{x}^i \mathbf{x} - 2\mathbf{a}^T \mathbf{x} + \mathbf{a}^T \mathbf{a} > R^2$$

使用核函数, 我们可以不限于球, 而是可以定义任意形状的边界。将点积用核函数替换, 我们得到(受限于相同的约束):

$$L_d = \sum_i \alpha^i K(\mathbf{x}^i, \mathbf{x}^i) - \sum_i \sum_s \alpha^i \alpha^s K(\mathbf{x}^i, \mathbf{x}^s) \quad (13.58)$$

例如, 使用二次多项式核可以使用任意的二次曲面。如果使用高斯核(参见式(13.30)), 则我们有局部球的并。我们把 \mathbf{x} 作为离群点而拒绝, 如果

$$K(\mathbf{x}, \mathbf{x}) - 2 \sum_i \alpha^i K(\mathbf{x}, \mathbf{x}^i) + \sum_i \sum_s \alpha^i \alpha^s K(\mathbf{x}^i, \mathbf{x}^s) > R^2$$

第三项不依赖于 \mathbf{x} , 因此它是常量(我们把它用作等量, 以便求解 R , 其中 \mathbf{x} 是边缘上的实例)。在高斯核的情况下, $K(\mathbf{x}, \mathbf{x}) = 1$, 该条件约简为对某个常数 R_c ,

$$\sum_i \alpha^i K_c(\mathbf{x}, \mathbf{x}^i) < R_c$$

除了解的稀疏性之外, 这类似于具有概率密度阈值 R_c 的核密度估计(参见 8.2.2 节, 参见图 13-11)。

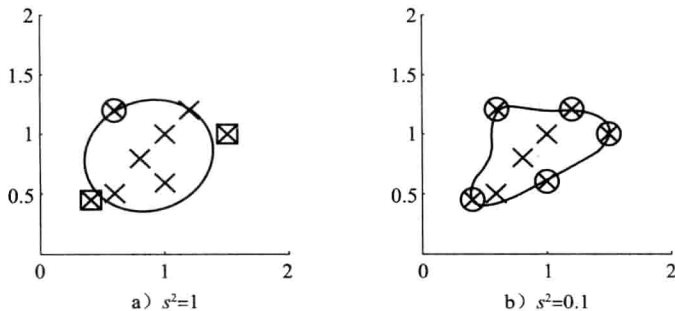


图 13-11 使用具有不同展宽的高斯核的一类支持向量机

一类支持向量机表示也有一种替代的、等价的 ν -SVM 类型, 它使用光滑的标准型 $\frac{1}{2} \|\mathbf{w}\|^2$ (Schölkopf 等 2001)。

13.12 核维度归约

由 6.3 节我们知道, 通过投影到协方差矩阵 Σ 的具有最大特征值的特征向量上, 主成分分析(PCA)降低维度。如果数据实例是中心化的($E[\mathbf{x}] = 0$), 则这可以记作 $\mathbf{X}^T \mathbf{X}$ 。在核版本中, 我们在 $\phi(\mathbf{x})$ 的而不是在原 \mathbf{x} 的空间处理, 因为通常这个新空间的维度 d 可能比数据集的大小 N 大得多, 我们宁愿使用 $N \times N$ 的矩阵 $\mathbf{X} \mathbf{X}^T$, 而不愿使用 $d \times d$ 的矩阵 $\mathbf{X}^T \mathbf{X}$ 。投影后的数据矩阵是 $\Phi = \phi(\mathbf{X})$, 因此我们在 $\Phi^T \Phi$ 特征向量上处理, 并进而在核矩阵 \mathbf{K} 的特征向量上处理。

核 PCA(kernel PCA)使用核矩阵的特征向量和特征值, 而这对应于在 $\phi(\mathbf{x})$ 空间中做线性维归约。当 \mathbf{c}_i 和 λ_i 对应于特征向量和特征值时, 投影后新的 k 维值可以用下式计算:

$$z_j^t = \sqrt{\lambda_j} \mathbf{c}_j^t, \quad j = 1, \dots, k, t = 1, \dots, N$$

图 13-12 给出了一个例子, 其中首先使用二次核, 然后使用核 PCA 把维度(从五维)降到二维, 并在那里实现线性 SVM。注意, 在一般情况下(例如使用高斯核), 特征值不一定衰减, 并且不能保证可以使用核 PCA 降低维度。

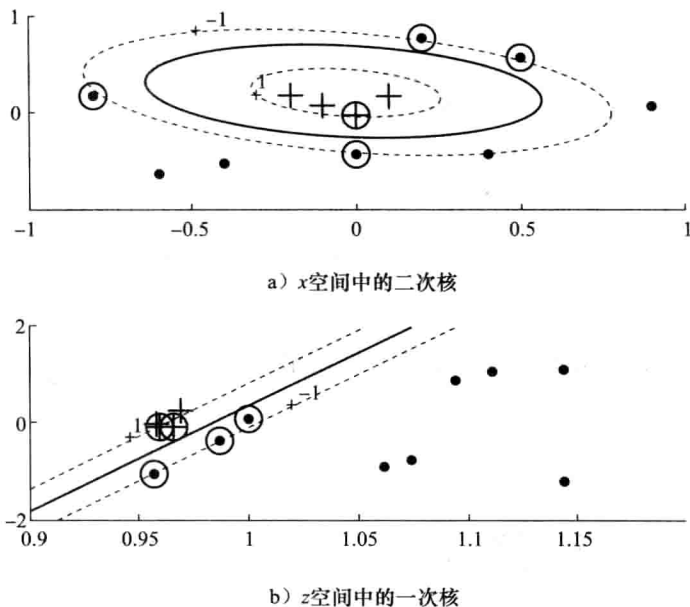


图 13-12 不使用原空间中的二次核(a), 而使用二次核值上的核 PCA 映射到二维新空间, 其中我们使用线性判别式(b); (五维中的)这两维贡献了方差的 80%

这里, 我们使用核值作为相似度值进行多维定标(参见 6.5 节)。例如, 取 $k = 2$, 我们可以在被核矩阵导出的空间中观察数据, 这可以让我们看到所使用的核定义的相似性效果如何。线性判别分析(LDA)也可以类似地核化(参见 6.6 节)(Müller 等 2001)。

第 6 章我们讨论了非线性维度归约方法 Isomap 和 LLE。事实上, 把式(6.47)中代价矩阵的元素看作输入对的核估计值, LLE 可以看作选择了特定核的核 PCA。当核函数定义为图

中测地距离时, 这一结论对 Isomap 也成立。

13.13 注释

通过非线性基函数把数据映射到新空间来推广线性模型的思想由来已久, 但是支持向量机的新颖性之处是把这一思想集成到学习算法中, 该学习算法的参数用数据实例子集定义(所谓对偶表示(dual representation)), 因此也不需要显式地计算基函数, 从而训练集的大小就限制了计算复杂度。对于高斯过程也是如此, 高斯过程中核函数称作协方差函数(参见 14.4 节)。

解的稀疏性表明与诸如 k -最近邻和 Parzen 窗口或高斯过程这样的非参数估计相比, 核方法更具优势, 并且使用核函数的灵活性使得我们能够处理非向量数据。由于优化问题存在唯一解, 因此不需要像神经网络那样使用迭代优化过程。由于所有这些理由, 支持向量机现在被看作最好的现成学习器, 并被广泛地用于许多领域, 特别是生物信息学(Schölkopf, Tsuda 和 Vert 2004)和自然语言处理应用, 在这些领域中正在开发越来越多的技巧来得到核函数(Shawe Taylor 和 Cristianini 2004)。

核函数的应用意味着不同的数据表示; 我们不再只是把实例(对象或事件)本身定义成属性的向量, 而是依据它与其他实例的相似或差异程度来定义它们; 这类似于使用距离矩阵(不必知道如何计算)的多维定标和使用某空间中向量的主成分分析之间的差别。

关于支持向量机的更多信息可以在 Vapnik(1995, 1998)、Schölkopf 和 Smola(2002)的书中找到。Cherkassky 和 Mulier 1998 中关于 SVM 的一章是非常值得一读的导论。Burges 1998、Smola 和 Schölkopf 1998 分别是 SVM 分类和回归的优秀指南。有一个专门的网站 <http://www.kernel-machines.org> 和许多免费软件包, 如 SVMlight(Joachims 2004)和 LIBSVM(Chang 和 Lin 2008)。

336
337

13.14 习题

1. 提出一种过滤算法, 找出非常可能不是支持向量的训练实例。
2. 在式(13.31)中, 如何估计 S ?
3. 在经验核映射中, 如何选择模板?
4. 在式(13.40)的局部化多核中, 为 $\eta_i(\mathbf{x} | \theta_i)$ 提出一个合适的模型, 并讨论如何训练它。
5. 在核回归中, ε 与噪声方差之间有何关系?
6. 在核回归中, 在偏倚和方差中使用不同的 ε 有什么影响?
7. 如何对分类使用一类 SVM?
8. 在如图 13-12 的情况中, 使用用高斯核的核 PCA。
9. 假定我们有相同对象的两种表示, 并且相互关联, 我们有不同的核。使用核 PCA, 如何使用这两种表示实现联合维度归约?

338

13.15 参考文献

- Allwein, E. L., R. E. Schapire, and Y. Singer. 2000. "Reducing Multiclass to Binary: A Unifying Approach for Margin Classifiers." *Journal of Machine Learning Research* 1: 113-141.
- Burges, C. J. C. 1998. "A Tutorial on Support Vector Machines for Pattern Recognition." *Data Mining and Knowledge Discovery* 2: 121-167.
- Chang, C.-C., and C.-J. Lin. 2008. *LIBSVM: A Library for Support Vector Machines*. <http://www.csie.ntu.edu.tw/~cjlin/libsvm/>.
- Cherkassky, V., and F. Mulier. 1998. *Learning from Data: Concepts, Theory, and Methods*. New York: Wiley.
- Cortes, C., and V. Vapnik. 1995. "Support Vector Networks." *Machine Learning* 20: 273-297.
- Dietterich, T. G., and G. Bakiri. 1995. "Solving Multiclass Learning Problems via Error-Correcting Output Codes." *Journal of Artificial Intelligence Research* 2: 263-286.
- Gönen, M., and E. Alpaydın. 2008. "Localized Multiple Kernel Learning." In *25th International Conference on Machine Learning*, ed. A. McCallum and S. Roweis, 352-359. Madison, WI: Omnipress.
- Jaakkola, T., and D. Haussler. 1999. "Exploiting Generative Models in Discriminative Classifiers." In *Advances in Neural Information Processing Systems 11*, ed. M. J. Kearns, S. A. Solla, and D. A. Cohn, 487-493. Cambridge, MA: MIT Press.
- Joachims, T. 2004. *SVMlight*, <http://svmlight.joachims.org>.
- Lanckriet, G. R. G., N. Cristianini, P. Bartlett, L. El Ghaoui, and M. I. Jordan. 2004. "Learning the Kernel Matrix with Semidefinite Programming." *Journal of Machine Learning Research* 5: 27-72.
- Mayoraz, E., and E. Alpaydın. 1999. "Support Vector Machines for Multiclass Classification." In *Foundations and Tools for Neural Modeling, Proceedings of IWANN'99, LNCS 1606*, ed. J. Mira and J. V. Sanchez-Andres, 833-842. Berlin: Springer.
- Müller, K. R., S. Mika, G. Rätsch, K. Tsuda, and B. Schölkopf. 2001. "An Introduction to Kernel-Based Learning Algorithms." *IEEE Transactions on Neural Networks* 12: 181-201.
- Noble, W. S. 2004. "Support Vector Machine Applications in Computational Biology." In *Kernel Methods in Computational Biology*, ed. B. Schölkopf, K. Tsuda, and J.-P. Vert, 71-92. Cambridge, MA: MIT Press.
- Platt, J. 1999. "Probabilities for Support Vector Machines." In *Advances in Large Margin Classifiers*, ed. A. J. Smola, P. Bartlett, B. Schölkopf, and D. Schuurmans, 61-74. Cambridge, MA: MIT Press.
- Schölkopf, B., J. Platt, J. Shawe-Taylor, A. J. Smola, and R. C. Williamson. 2001. "Estimating the Support of a High-Dimensional Distribution." *Neural Computation* 13: 1443-1471.
- Schölkopf, B., and A. J. Smola. 2002. *Learning with Kernels: Support Vector Machines, Regularization, Optimization, and Beyond*. Cambridge, MA: MIT Press.

- Schölkopf, B., A. J. Smola, R. C. Williamson, and P. L. Bartlett. 2000. "New Support Vector Algorithms." *Neural Computation* 12: 1207-1245.
- Schölkopf, B., K. Tsuda, and J.-P. Vert, eds. 2004. *Kernel Methods in Computational Biology*. Cambridge, MA: MIT Press.
- Shawe-Taylor, J., and N. Cristianini. 2004. *Kernel Methods for Pattern Analysis*. Cambridge, UK: Cambridge University Press.
- Smola, A., and B. Schölkopf. 1998. *A Tutorial on Support Vector Regression*, NeuroCOLT TR-1998-030, Royal Holloway College, University of London, UK.
- Sonnenburg, S., G. Rätsch, C. Schäfer, and B. Schölkopf. 2006. "Large Scale Multiple Kernel Learning." *Journal of Machine Learning Research* 7: 1531-1565.
- Tax, D. M. J., and R. P. W. Duin. 1999. "Support Vector Domain Description." *Pattern Recognition Letters* 20: 1191-1199.
- Vapnik, V. 1995. *The Nature of Statistical Learning Theory*. New York: Springer.
- Vapnik, V. 1998. *Statistical Learning Theory*. New York: Wiley.
- Vert, J.-P., K. Tsuda, and B. Schölkopf. 2004. "A Primer on Kernel Methods." In *Kernel Methods in Computational Biology*, ed. B. Schölkopf, K. Tsuda, and J.-P. Vert, 35-70. Cambridge, MA: MIT Press.
- Weston, J., and C. Watkins. 1998. "Multiclass Support Vector Machines." *Technical Report CSD-TR-98-04*, Department of Computer Science, Royal Holloway, University of London.

贝叶斯估计

在贝叶斯方法中，我们把参数看作具有先验分布的随机变量。我们继续 4.4 节中断的讨论，并讨论 3 种情况：估计分布的参数、估计模型的参数和高斯过程。

14.1 引言

当我们有关于参数的某些先验信息时，可以使用贝叶斯估计 (Bayesian estimation)。例如，在观察样本估计分布的均值 μ 之前，我们可能有先验知识，知道均值接近于 2，在 1 和 3 之间。当我们只有小样本时，这种先验知识特别重要。在这种情况下，我们感兴趣的是把数据告诉我们的（即由样本计算的值）与我们的先验信息结合。

4.2 节讨论的最大似然方法将参数看作未知的常量。在贝叶斯估计中，正如我们在 4.4 节开始讨论的，参数被看作一个随机变量，使得我们可以使用先验概率分布 (prior probability distribution) 对任何先验信息编码。例如，知道 μ 很可能在 1 和 3 之间，我们可以这样记 $p(\mu)$ ，其密度的大部分落在区间 $[1, 3]$ 。

使用贝叶斯规则，我们把先验与似然结合在一起，并计算后验概率分布 (posterior probability distribution)：

$$p(\theta | \mathcal{X}) = \frac{p(\theta)p(\mathcal{X} | \theta)}{p(\mathcal{X})} \quad (14.1)$$

$p(\theta)$ 是先验密度，它是我们在考察样本之前知道的 θ 可能的取值。 $p(\mathcal{X} | \theta)$ 是样本似然 (sample likelihood)，它告诉我们如果分布的参数取该 θ 值，样本 \mathcal{X} 有多大可能出现。例如，如果样本中的实例落在 5 与 10 之间，那么若 μ 为 7 则这种样本是可能的，但若 μ 为 3 则不大可能，而 μ 为 1 则更不可能。分母中的 $p(\mathcal{X})$ 是正规化子，以确保后验 (posterior) $p(\theta | \mathcal{X})$ 的积分等于 1。 $p(\theta | \mathcal{X})$ 又称后验概率，因为它告诉我们在看到样本之后 θ 取特定值的可能性有多大。贝叶斯规则取先验分布，把它与数据揭示的结合在一起，并产生后验分布。然后，我们在稍后的推断中使用这个后验分布。

例如，假设我们有从具有某个未知参数 θ 的分布中提取的旧样本 \mathcal{X} 。然后，我们可以再抽取一个实例 x ，并计算它的概率分布。我们可以把这形象地表示为一个图模型 (参见第 16 章)，如图 14-1 所示，描绘的是一个生成模型 (generative model)，它表示数据如何生成：首先由 $p(\theta)$ 选

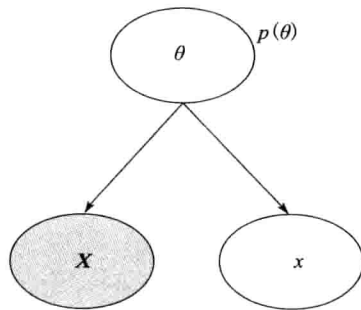


图 14-1 这是一个生成图模型。弧是抽样方向：首先由 $p(\theta)$ 选 θ ，然后通过从 $p(x' | \theta)$ 中抽样生成数据。给定 θ ，新实例 x 和过去的样本 \mathcal{X} 是独立的：这是独立同分布假设。如果我们不知道 θ ，则它们是依赖的：我们使用贝叶斯规则由给定的 \mathcal{X} (用阴影显示) 推断 θ ，这逆转了方向以便计算 $p(\theta | \mathcal{X})$ ，然后用它来填充 x

θ , 然后把它用于样本 \mathcal{X} , 再用于新实例 x 。我们把联合概率表示成

$$p(x, \mathcal{X}, \theta) = p(\theta)p(\mathcal{X}|\theta)p(x|\theta)$$

给定以往的样本 \mathcal{X} , 我们用它估计新实例 x 的概率:

$$\begin{aligned} p(x|\mathcal{X}) &= \frac{p(x, \mathcal{X})}{p(\mathcal{X})} = \frac{\int_{\theta} p(x, \mathcal{X}, \theta) d\theta}{p(\mathcal{X})} = \frac{\int_{\theta} p(\theta)p(\mathcal{X}|\theta)p(x|\theta) d\theta}{p(\mathcal{X})} \\ &= \int_{\theta} p(\theta|\mathcal{X})p(x|\theta) d\theta \end{aligned} \quad (14.2)$$

在计算 $p(\theta|\mathcal{X})$ 时, 贝叶斯规则允许我们逆转弧的方向并做诊断推理。然后, 推断出的 θ 的分布用来导出新 x 的预测分布。

我们看到, 我们的估计是 θ 所有可能值估计的加权和(如果 θ 是离散值, 则我们用 \sum_{θ} 替换 $\int d\theta$), 权重是给定 \mathcal{X} , θ 的可能性。

这是完全贝叶斯处理(full Bayesian treatment); 如果后验概率不容易求积分, 则这样做或许不可能。正如我们在 4.4 节看到的, 在最大化后验(Maximum A Posteriori, MAP)估计中, 我们使用后验的众数:

$$\theta_{\text{MAP}} = \arg \max_{\theta} p(\theta|\mathcal{X}) \text{ 并且 } p_{\text{MAP}}(x|\mathcal{X}) = p(x|\theta_{\text{MAP}})$$

MAP 估计对应于假定后验在单点(即众数)附近制造很窄的尖峰。如果先验 $p(\theta)$ 在所有的 θ 上是均匀的, 则后验 $p(\theta|\mathcal{X})$ 的众数与似然 $p(\mathcal{X}|\theta)$ 的众数在同一点上, 并且 MAP 估计与最大似然(maximum likelihood, ML)估计等价。这意味使用 ML 对应于假定 θ 的不同值之间没有先验分布。

现在, 让我们看看贝叶斯估计如何用于不同的分布和应用类型。

14.2 分布参数的估计

14.2.1 离散变量

假定每个实例都是一个多元变量, 取 K 个不同状态之一(参见 4.2.2 节)。我们说 $x_i^t = 1$, 如果实例 t 在状态 i 中, 并且对于任意 $j \neq i$, $x_j^t = 0$ 。参数是状态的概率 $\mathbf{q} = [q_1, q_2, \dots, q_K]^T$, 其中 $q_i (i=1, \dots, K)$ 满足 $q_i \geq 0$, 并且 $\sum_i q_i = 1$ 。

样本似乎是

$$p(\mathcal{X}|\mathbf{q}) = \prod_{t=1}^N \prod_{i=1}^K q_i^{x_i^t}$$

我们使用的先验分布是狄利克雷分布(Dirichlet distribution)

$$\text{Dirichlet}(\mathbf{q}|\boldsymbol{\alpha}) = \frac{\Gamma(\alpha_0)}{\Gamma(\alpha_1) \dots \Gamma(\alpha_K)} \prod_{i=1}^K q_i^{\alpha_i-1}$$

其中, $\boldsymbol{\alpha} = [\alpha_1, \dots, \alpha_K]^T$, 并且 $\alpha_0 = \sum_i \alpha_i$ 。 α_i 是先验参数, 称作超参数(hyperparameter)。 $\Gamma(x)$ 是伽马函数(Gamma function), 定义为

$$\Gamma(x) \equiv \int_0^{\infty} u^{x-1} e^{-u} du$$

例如, x^t 可以对应于新闻文档, 状态可以对应于 K 个不同的新闻类别: 体育、政治、艺术等。于是, 概率 q_i 对应于不同类别新闻所占的比例, 而关于它们的先验使得我们可以对这些比例的先验编码。例如, 我们可以预期与体育相关的新闻比与艺术相关的新闻更多。

给定先验和似然, 我们可以导出后验

$$p(\mathbf{q} | \mathcal{X}) \propto p(\mathcal{X} | \mathbf{q}) p(\mathbf{q} | \boldsymbol{\alpha}) \propto \prod_i q_i^{\alpha_i + N_i - 1} \quad (14.3)$$

其中, $N_i = \sum_{t=1}^N x_i^t$ 。我们看到后验具有与先验相同的形式, 并且我们称这种先验为共轭先验 (conjugate prior)。先验和似然都是 q_i 幂的乘积形式, 并且我们可以把它们结合以构成后验

$$p(\mathbf{q} | \mathcal{X}) = \frac{\Gamma(\alpha_0 + N)}{\Gamma(\alpha_1 + N_1) \dots \Gamma(\alpha_K + N_K)} \prod_{i=1}^K q_i^{\alpha_i + N_i - 1} = \text{Dirichlet}(\mathbf{q} | \boldsymbol{\alpha} + \mathbf{n}) \quad (14.4)$$

其中, $\mathbf{n} = [N_1, \dots, N_K]^T$, $\sum_i N_i = N$ 。

观察式(14.3)我们可以得到超参数 α_i 的一种解释(Bishop 2006)。正如 n_i 是 N 个样本中状态 i 出现的次数一样, 我们可以将 α_i 看作在 α_0 个实例的某个假想样本中状态 i 出现的次数。在定义先验时我们主观地认为: 在 α_0 的样本中, 我们期望它们之中的 α_i 个属于状态 i 。注意, 较大的 α_0 意味着对我们的主观比例有较高的置信度(更尖的分布): 我们期望 100 次出现中的 60 次属于状态 1 的置信度比期望 10 次出现中的 6 次属于状态 1 的置信度高。于是, 后验是另一个狄利克雷, 它对分别由先验和似然给定的想象和实际状态出现的次数求和。

共轭性具有很好的含义。在顺序情况下, 我们接收实例序列。因为后验与先验具有相同的形式, 因此当前后验由过去的所有实例累积信息, 并且成为下一个实例的先验。

当变量是二元时, $x^t \in \{0, 1\}$, 多元样本变成贝努利

$$p(\mathcal{X} | q) = \prod_i q^{x^i} (1 - q)^{1 - x^i}$$

并且狄利克雷先验归约为贝塔分布(beta distribution)

$$\text{beta}(q | \alpha, \beta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} q^{\alpha-1} (1 - q)^{\beta-1}$$

例如, x^t 可以为 0 或 1, 取值依赖于大小为 N 的随机样本中索引为 t 的电子邮件是正常邮件还是垃圾邮件。于是, 定义 q 上的先验使得我们可以对垃圾邮件出现的概率定义先验信念: 在平均情况下, 我们预料电子邮件中有 $\alpha/(\alpha + \beta)$ 封是垃圾邮件。

贝塔是共轭先验, 并且对于该后验我们得到

$$p(q | A, N, \alpha, \beta) \propto q^{A + \alpha - 1} (1 - q)^{N - A + \beta - 1}$$

其中 $A = \sum_t x^t$, 并且再次看到我们组合了想象和实际样本中的出现。注意, 当 $\alpha = \beta = 1$ 时, 有均匀先验, 并且后验与似然有相同的形状。随着这两个计数(无论是关于先验的 α 和 β , 还是关于后验的 $\alpha + A$ 和 $\beta + N - A$) 的增加和它们之差的增加, 我们得到具有更小方差的更尖分布(参见图 14-2)。随着我们看到更多

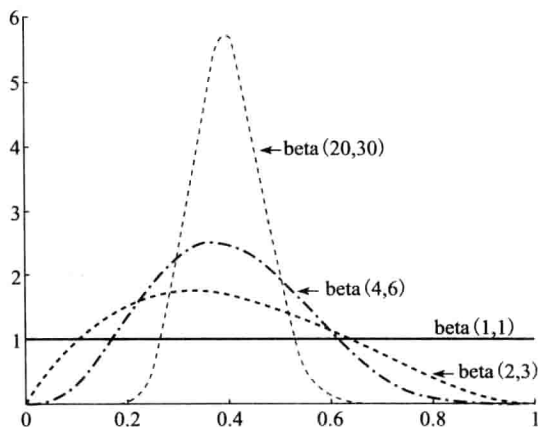


图 14-2 不同 (α, β) 的贝塔分布图

的数据(想象的或实际的), 方差减小。

14.2.2 连续变量

现在我们考虑实例是高斯分布 $p(x) \sim \mathcal{N}(\mu, \sigma^2)$, 参数是 μ 和 σ^2 的情况; 我们已经在 4.4 节简略讨论过。样本似然为

$$p(\mathcal{X} | \mu, \sigma^2) = \prod_i \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{(x^i - \mu)^2}{2\sigma^2}\right] \quad (14.5)$$

μ 的共轭先验是高斯 $p(\mu) \sim \mathcal{N}(\mu_0, \sigma_0^2)$, 并且后验写作

$$p(\mu | \mathcal{X}) \propto p(\mu)p(\mathcal{X} | \mu) \sim \mathcal{N}(\mu_N, \sigma_N^2)$$

其中

$$\mu_N = \frac{\sigma^2}{N\sigma_0^2 + \sigma^2}\mu_0 + \frac{N\sigma^2}{N\sigma_0^2 + \sigma^2}m \quad (14.6)$$

$$\frac{1}{\sigma_N^2} = \frac{1}{\sigma_0^2} + \frac{N}{\sigma^2} \quad (14.7)$$

其中 $m = \sum_i x^i / N$ 是样本平均值。我们看到后验密度的均值(贝叶斯估计) μ_N , 是先验均值 μ_0 和样本均值 m 的加权平均, 其中权重与它们的方差成反比(例子参见图 14-3)。注意, 因为两个系数都在 0 和 1 之间并且其和为 1, 所以 μ_N 总是在 μ_0 和 m 之间。当样本规模 N 或先验 σ_0^2 的方差大的时候, 贝叶斯估计接近于 m , 更多地依赖样本提供的信息。当 σ_0^2 小时, 即当我们对 μ 正确值先验的不确定性较小时, 或当我们有小样本时, 我们的先验猜测 μ_0 具有更大影响。

当 σ_0 或 σ 变小, 或 N 较大时, σ_N 变小。还要注意, σ_N 小于 σ_0 和 σ/\sqrt{N} , 即后验方差小于先验方差和 m 的方差。将这两个结果结合在后验估计中比单独使用先验或样本更好。

对于方差, 我们处理的精度(precision)是方差的倒数, $\lambda \equiv 1/\sigma^2$ 。使用这一事实, 样本似然写成

$$\begin{aligned} p(\mathcal{X} | \lambda) &= \prod_i \frac{\lambda^{1/2}}{\sqrt{2\pi}} \exp\left[-\frac{\lambda}{2}(x^i - \mu)^2\right] \\ &= \lambda^{N/2} (2\pi)^{-N/2} \exp\left[-\frac{\lambda}{2} \sum_i (x^i - \mu)^2\right] \end{aligned} \quad (14.8)$$

精度的共轭先验是伽马分布(Gamma distribution):

$$p(\lambda) \sim \text{Gamma}(a_0, b_0) = \frac{1}{\Gamma(a_0)} b_0^{a_0} \lambda^{a_0-1} \exp(-b_0 \lambda)$$

而后验是

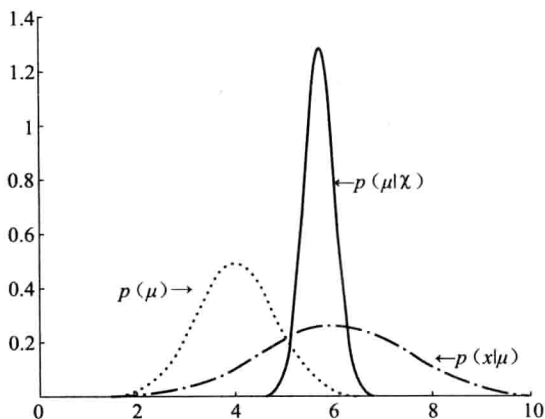


图 14-3 20 个数据点从 $p(x) \sim \mathcal{N}(6, 1.5^2)$ 中抽取, 先验是 $p(\mu) \sim \mathcal{N}(4, 0.8^2)$, 于是后验是 $p(\mu | \mathcal{X}) \sim \mathcal{N}(5.7, 0.3^2)$

$$p(\lambda | \chi) \propto p(\chi | \lambda) p(\lambda) \sim \text{Gamma}(a_N, b_N)$$

其中

$$a_N = a_0 + N/2 \quad (14.9)$$

$$b_N = b_0 + \frac{N}{2} s^2$$

其中 $s^2 = \sum_i (x^i - \mu)^2 / N$ 是样本方差。我们再一次看到后验估计是先验和样本统计量的加权和。

14.3 函数参数的贝叶斯估计

现在, 我们对回归和分类讨论参数估计, 不讨论分布的参数, 而是讨论输入的某个函数的参数。我们的方法仍然是将这些参数看作具有一种先验分布的随机变量, 并使用贝叶斯规则计算后验分布。然后, 或者求积分、近似它, 或者使用 MAP 估计。

14.3.1 回归

让我们考虑线性回归模型的情况

$$\mathbf{r} = \mathbf{w}^T \mathbf{x} + \varepsilon, \text{ 其中 } \varepsilon \sim \mathcal{N}(0, \beta^{-1}) \quad (14.10)$$

其中 β 是加法噪声的精度。

参数是权重 \mathbf{w} , 并且我们有样本 $\chi = \{\mathbf{x}^i, \mathbf{r}^i\}_{i=1}^N$, 其中 $\mathbf{x} \in \mathcal{R}^d$, $\mathbf{r}^i \in \mathcal{R}$ 。我们可以把 χ 分解成输入矩阵和期望输出的向量 $\chi = [\mathbf{X}, \mathbf{r}]$ 。由式(14.10), 我们有

$$p(\mathbf{r}^i | \mathbf{x}^i, \mathbf{w}, \beta) \sim \mathcal{N}(\mathbf{w}^T \mathbf{x}, \beta^{-1})$$

前面, 在 4.6 节我们看到对数似然是

$$\begin{aligned} \mathcal{L}(\chi | \mathbf{w}) &\equiv \log p(\chi | \mathbf{w}) = \log p(\mathbf{r}, \mathbf{X} | \mathbf{w}) \\ &= \log p(\mathbf{r} | \mathbf{X}, \mathbf{w}) + \log p(\mathbf{X}) \end{aligned}$$

其中第二项是常数, 独立于参数。我们把第一项展开成

$$\begin{aligned} \mathcal{L}(\mathbf{r} | \mathbf{X}, \mathbf{w}, \beta) &= \log \prod_i p(\mathbf{r}^i | \mathbf{x}^i, \mathbf{w}, \beta) \\ &= -N \log(\sqrt{2\pi}) + N \log \beta - \frac{\beta}{2} \sum_i (\mathbf{r}^i - \mathbf{w}^T \mathbf{x}^i)^2 \end{aligned} \quad (14.11)$$

对于 ML 估计, 我们找出最大化上式的 \mathbf{w} , 或等价地, 最小化上式的最后一项, 即误差的平方和, 该项可以改写为

$$E = (\mathbf{r} - \mathbf{X}\mathbf{w})^T (\mathbf{r} - \mathbf{X}\mathbf{w})$$

关于 \mathbf{w} 求导并令它等于 0, 我们得到最大似然估计(5.8 节曾经推导出它):

$$\mathbf{w}_{ML} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{r} \quad (14.12)$$

计算出参数之后我们就可以做预测。给定新的输入 \mathbf{x}' , 计算得到的响应为

$$\mathbf{r}' = \mathbf{w}_{ML}^T \mathbf{x}' \quad (14.13)$$

对于非线性模型 $g(\mathbf{x} | \mathbf{w})$, 例如对于多层感知器, 其中 \mathbf{w} 是权重, 我们使用梯度下降最小化

$$E(\chi | \mathbf{w}) = [\mathbf{r}^i - g(\mathbf{x}^i | \mathbf{w})]^2$$

并且将最小化上式的 \mathbf{w}_{LSQ} 称作最小二乘估计子 (least square estimator)。于是, 预测用下式计算:

$$\mathbf{r}' = g(\mathbf{x}' | \mathbf{w}_{LSQ})$$

在贝叶斯方法的情况下, 对于参数, 我们定义一个高斯先验 (Gaussian prior):

$$p(\mathbf{w}) \sim \mathcal{N}(\mathbf{0}, \alpha^{-1} \mathbf{I})$$

它是共轭先验, 并且对于后验, 我们得到

$$p(\mathbf{w} | \mathcal{X}) \sim \mathcal{N}(\boldsymbol{\mu}_N, \boldsymbol{\Sigma}_N)$$

其中

$$\boldsymbol{\mu}_N = \beta \boldsymbol{\Sigma}_N \mathbf{X}^T \mathbf{r} \quad (14.14)$$

$$\boldsymbol{\Sigma}_N = (\alpha \mathbf{I} + \beta \mathbf{X}^T \mathbf{X})^{-1} \quad (14.15)$$

为了计算总输出, 我们在后验上积分

$$\mathbf{r}' = \int \mathbf{w}^T \mathbf{x}' p(\mathbf{w} | \mathcal{X}) d\mathbf{w}$$

如果我们想用点估计, 则 MAP (或贝叶斯, 因为后验是高斯的) 估计是

$$\mathbf{w}_{MAP} = \boldsymbol{\mu}_N = \beta (\alpha \mathbf{I} + \beta \mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{r} \quad (14.16)$$

并且我们将密度替换成单点, 即均值

$$\mathbf{r}' = \mathbf{w}_{MAP}^T \mathbf{x}'$$

方差为

$$\text{Var}(\mathbf{r}') = \beta^{-1} + (\mathbf{x}')^T \boldsymbol{\Sigma}_N \mathbf{x}' \quad (14.17)$$

将式 (14.16) 与式 (14.14) 的 ML 估计进行比较, 这可以看作正则化——我们在对角线增加了一个常数项 α , 以便使得矩阵可逆。

先验 $p(\mathbf{w}) \sim \mathcal{N}(\mathbf{0}, \alpha^{-1} \mathbf{I})$ 表明我们期望参数接近于 0, 展宽与 α 成反比。当 $\alpha \rightarrow 0$ 时, 我们有平坦的先验, 并且 MAP 估计收敛于 ML 估计。

在图 14-4 中我们看到, 如果增大 α , 则迫使参数更接近 0, 并且后验分布移近原点并收缩。如果减小 β , 则假定噪声具有高方差, 并且后验也具有高方差。

如果我们取后验的对数, 则有

$$\begin{aligned} \log p(\mathbf{w} | \mathbf{X}, \mathbf{r}) &\propto \log p(\mathbf{X}, \mathbf{r} | \mathbf{w}) + \log p(\mathbf{w}) \\ &\propto \log p(\mathbf{r} | \mathbf{X}, \mathbf{w}) + \log p(\mathbf{w}) \\ &= -\frac{\beta}{2} \sum_i (\mathbf{r}'^i - \mathbf{w}^T \mathbf{x}'^i)^2 - \frac{\alpha}{2} \mathbf{w}^T \mathbf{w} + c \end{aligned}$$

我们对其最大化, 以得到 MAP 估计。在一般情况下, 给定模型 $g(\mathbf{x} | \mathbf{w})$, 我们可以写一个增广的误差函数

$$E_{\text{ridge}}(\mathbf{w} | \mathcal{X}) = \sum_i [\mathbf{r}'^i - g(\mathbf{x}'^i | \mathbf{w})]^2 + \lambda \sum_i w_i^2$$

其中 $\lambda \equiv \alpha/\beta$ 。在统计学中, 这称作参数收缩 (parameter shrinkage) 或岭回归 (ridge regression)。在 4.8 节, 我们称为正则化 (regularization), 而在 11.9 节, 我们称这为神经网络中的权衰减 (weight decay)。第一项是似然的负对数, 而第二项是处罚远离 0 的 w_i (正如先验的 α 所示)。

尽管这种方法减少 $\sum_i w_i^2$, 但是它并不强制个体 w_i 为 0; 即它不能用于特征选择, 也就是说, 不能用于决定哪些 x_i 是冗余的。为此, 可以使用利用 L_1 范数而不是 L_2 范数的拉普拉

斯先验 (Laplacian prior) (Figueiredo 2003):

$$p(\mathbf{w} | \alpha) = \prod_i \frac{\alpha}{2} \exp(-\alpha |w_i|) = \left(\frac{\alpha}{2}\right)^d \exp(-\alpha \sum_i |w_i|)$$

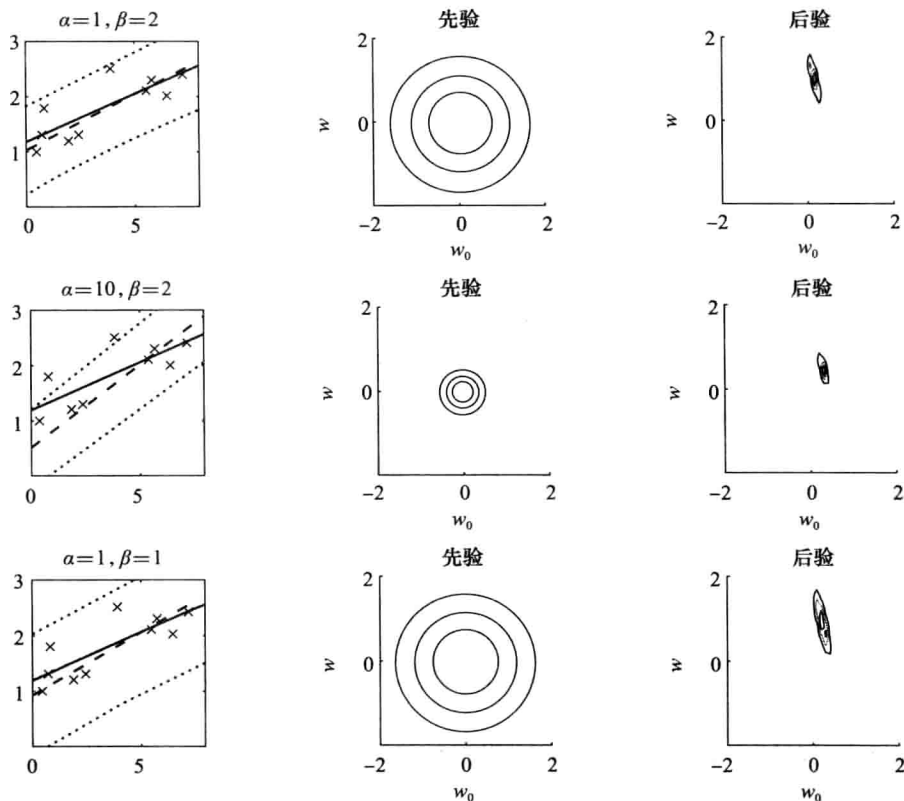


图 14-4 不同 α 和 β 值的贝叶斯线性回归。左边：“ \times ”是数据点，直线是 ML 解，还用虚线显示了具有一个标准差的误差条线的 MAP 解。中间：中心在 0、方差为 $1/\alpha$ 的先验密度。右边：后验密度，其均值是 MAP 解。我们看到，当增大 α 时，先验的方差收缩，并且线移向平坦的 0 线。当 β 减小时，假定有更多的噪声，并且后验密度具有较高方差

后验概率不再是高斯的，而 MAP 估计通过最小化下式找出：

$$E_{\text{lasso}}(\mathbf{w} | \mathcal{X}) = \sum_i (r^i - \mathbf{w}^T \mathbf{x}^i)^2 + 2\sigma^2 \alpha \sum_i |w_i|$$

其中 σ^2 是噪声方差 (对此我们插入我们的估计)。这称作 lasso (least absolute shrinkage and selection operator) (Tibshirani 1996)。为了明白为什么 L_1 降低了稀疏性，我们考虑两个权重 $[w_1, w_2]^T$ 的情况 (Figueiredo 2003)： $\|[1, 0]^T\|_2 = \|[1/\sqrt{2}, 1/\sqrt{2}]^T\|_2 = 1$ 而 $\|[1, 0]^T\|_1 = 1 < \|[1/\sqrt{2}, 1/\sqrt{2}]^T\|_1 = \sqrt{2}$ ，因此 L_1 更倾向于置 w_2 为 0 并使用较大的 w_1 ，而不是让它们都取较小的值。

14.3.2 基函数或核函数的使用

使用式 (14.14) 的贝叶斯估计，预测可以表示为

$$\begin{aligned}
 r' &= (\mathbf{x}')^T \mathbf{w} \\
 &= \beta(\mathbf{x}')^T \Sigma_N \mathbf{X}^T \mathbf{r} \\
 &= \sum_i \beta(\mathbf{x}')^T \Sigma_N \mathbf{x}^i r^i
 \end{aligned}$$

这是对偶表示(dual representation)。当我们可以用训练数据,或像支持向量机(参见第13章)那样用训练数据的一个子集表示参数时,我们可以把预测写成当前输入和过去数据的函数。我们可以表示为

$$r' = \sum_i K(\mathbf{x}', \mathbf{x}^i) r^i \quad (14.18)$$

其中,我们定义

$$K(\mathbf{x}', \mathbf{x}^i) = \beta(\mathbf{x}')^T \Sigma_N \mathbf{x}^i \quad (14.19)$$

我们可以推广式(14.19)中的线性核,使用非线性基函数 $\phi(\mathbf{x})$ 映射到新空间,在新空间中我们拟合线性模型。在这种情况下,我们有 k 维 $\phi(\mathbf{x})$ 而不是 d 维 \mathbf{x} , 其中 k 是基函数的个数,并且我们有 $N \times k$ 的基函数的像 Φ , 而不是 $N \times d$ 的数据矩阵 \mathbf{X} 。

在检验期间,我们有

$$\begin{aligned}
 r' &= \phi(\mathbf{x}')^T \mathbf{w}, \text{ 其中 } \mathbf{w} = \beta \Sigma_N^\phi \Phi^T \mathbf{r}, \Sigma_N^\phi = (\alpha \mathbf{I} + \beta \Phi^T \Phi)^{-1} \\
 &= \beta \phi(\mathbf{x}')^T \Sigma_N^\phi \Phi^T \mathbf{r} \\
 &= \sum_i \beta \phi(\mathbf{x}')^T \Sigma_N^\phi \phi(\mathbf{x}^i) r^i \\
 &= \sum_i K(\mathbf{x}', \mathbf{x}^i) r^i
 \end{aligned} \quad (14.20)$$

其中,我们定义

$$K(\mathbf{x}', \mathbf{x}^i) = \beta \phi(\mathbf{x}')^T \Sigma_N^\phi \phi(\mathbf{x}^i) \quad (14.21)$$

作为等价核。这是 $\phi(\mathbf{x})$ 的空间中的对偶表示。我们可以将估计表示成训练集中实例影响的加权和,其中影响由核函数 $K(\mathbf{x}', \mathbf{x}^i)$ 给定;这类似于我们在第8章讨论的非参数的核光滑,或第13章的核机器。

误差条线可以用下式定义

$$\text{Var}(r') = \beta^{-1} + \phi(\mathbf{x}')^T \Sigma_N^\phi \phi(\mathbf{x}')$$

对于线性、二次核、四次核,图14-5给出了一个例子。

正如在回归中我们可以在原 \mathbf{x} 或 $\phi(\mathbf{x})$ 上处理一样,在贝叶斯回归中,我们也可以在预处理过的 $\phi(\mathbf{x})$ 上处理,定义该空间中的参数。本章稍后我们将考察高斯过程,在该过程中我们可以直接定义和使用 $K(\mathbf{x}, \mathbf{x}')$, 而不必计算 $\phi(\mathbf{x})$ 。

14.3.3 贝叶斯分类

在两类问题中,我们有单输出,并且假定为线性模型,则有

$$P(C_1 | \mathbf{x}^i) = y^i = \text{sigmoid}(\mathbf{w}^T \mathbf{x}^i)$$

贝努利样本的对数似然为

$$\mathcal{L}(\mathbf{r} | \mathbf{X}) = \sum_i r^i \log y^i + (1 - r^i) \log (1 - y^i)$$

例如,我们使用梯度下降最大化它,或最小化它的负对数(互熵),得到 ML 估计。这称作逻辑斯谛判别式(logistic discrimination)(参见10.7节)。

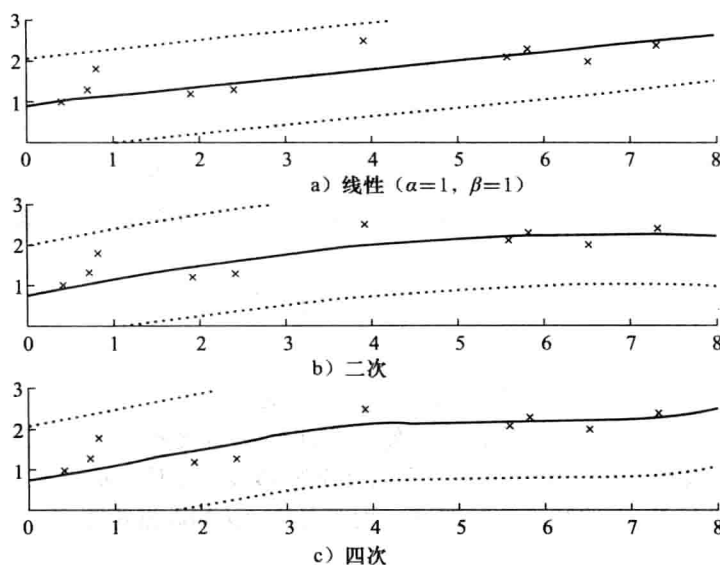


图 14-5 使用具有一个标准差误差条线核的贝叶斯回归: (a) 线性: $\phi(\mathbf{x}) = [1, x]^T$; (b) 二次: $\phi(\mathbf{x}) = [1, x, x^2]^T$; (c) 四次: $\phi(\mathbf{x}) = [1, x, x^2, x^3, x^4]^T$

在贝叶斯方法中, 我们假定高斯先验

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{m}_0, \mathbf{S}_0) \quad (14.22)$$

而后验的对数为

$$\begin{aligned} \log p(\mathbf{w} | \mathbf{r}, \mathbf{X}) &\propto \log p(\mathbf{w}) + \log p(\mathbf{r} | \mathbf{w}, \mathbf{X}) \\ &= -\frac{1}{2}(\mathbf{w} - \mathbf{m}_0)^T \mathbf{S}_0^{-1}(\mathbf{w} - \mathbf{m}_0) \\ &\quad + \sum_i r^i \log y^i + (1 - r^i) \log (1 - y^i) + c \end{aligned} \quad (14.23)$$

这个后验分布不再是高斯的, 并且我们不能精确地求积分。我们可以使用拉普拉斯近似 (Laplace approximation), 方法如下 (MacKay 2003)。假设我们想要近似某个分布 $f(x)$, 不必是规范化的 (积分为 1)。在拉普拉斯近似中, 我们找出 $f(x)$ 的众数 x_0 , 拟合一个中心在 x_0 的高斯函数 $q(x)$, 而后如果我们想要积分, 就在拟合的高斯函数上积分。为了得到该高斯的方差, 我们考虑 $f(\cdot)$ 在 $x = x_0$ 处的泰勒展开

$$\log f(x) = \log f(x_0) - \frac{1}{2}a(x - x_0)^2 + \dots$$

其中

$$a \equiv -\frac{d}{dx^2} \log f(x) \Big|_{x=x_0}$$

注意, 第一项 (线性项) 消失, 因为在众数上的一阶导数为 0。取指数我们得到

$$f(x) = f(x_0) \exp \left[-\frac{a}{2}(x - x_0)^2 \right]$$

为了规范化 $f(x)$, 我们考虑高斯分布

$$\int \frac{1}{\sqrt{2\pi}(1/\sqrt{a})} \exp \left[-\frac{a}{2}(x - x_0)^2 \right] dx = 1 \Rightarrow \int \exp \left[-\frac{a}{2}(x - x_0)^2 \right] dx = \sqrt{a/2\pi}$$

因此

$$q(x) = \sqrt{a/2\pi} \exp\left[-\frac{a}{2}(x-x_0)^2\right] \sim \mathcal{N}(x_0, 1/a)$$

在多元情况下, $x \in \mathbb{R}^d$, 我们有

$$\log f(x) = \log f(x_0) - \frac{1}{2}(x-x_0)^T \mathbf{A}(x-x_0) + \dots$$

其中 \mathbf{A} 是二阶导数的 (Hessian) 矩阵:

$$\mathbf{A} = -\nabla \nabla \log f(x) \big|_{x=x_0}$$

于是, 拉普拉斯近似为

$$f(x) = \frac{|\mathbf{A}|^{1/2}}{(2\pi)^{d/2}} \exp\left[-\frac{1}{2}(x-x_0)^T \mathbf{A}(x-x_0)\right] \sim \mathcal{N}_d(x_0, \mathbf{A}^{-1})$$

在讨论了如何近似之后, 现在可以使用它计算后验密度。 \mathbf{w}_{MAP} 是 $p(\mathbf{w} | \mathbf{r}, \mathbf{X})$ 的众数, 取作均值, 而协方差矩阵由负的对数似然的二阶导数矩阵的逆给出:

$$\mathbf{S}_N = -\nabla \nabla \log p(\mathbf{w} | \mathbf{r}, \mathbf{X}) = \mathbf{S}_0^{-1} + \sum_i y'(1-y') \mathbf{x}'(\mathbf{x}')^T$$

于是, 我们在这个高斯函数上积分, 以估计类概率

$$P(C_1 | \mathbf{x}) = y = \int \text{sigmoid}(\mathbf{w}^T \mathbf{x}) q(\mathbf{w}) d\mathbf{w}$$

其中 $q(\mathbf{w}) \sim \mathcal{N}(\mathbf{w}_{MAP}, \mathbf{S}_N^{-1})$ 。另一个难题是我们不能解析地求解带 sigmoid 的高斯卷积的积分。概率单位函数 (probit function) 与 sigmoid 函数具有相同的 S 形, 如果我们代之以概率单位函数, 则可以得到解析解 (Bishop 2006)。

354
355

14.4 高斯过程

假定我们有线性模型 $y = \mathbf{w}^T \mathbf{x}$ 。于是, 对于每个 \mathbf{w} , 我们有一条直线。给定先验分布 $p(\mathbf{w})$, 我们得到一个直线分布, 或更具体地说, 对于任意的 \mathbf{w} , 当 \mathbf{w} 是从 $p(\mathbf{w})$ 中抽样时, 我们得到在 \mathbf{x} 处计算的 y 值 (记作 $y(\mathbf{x} | \mathbf{w})$) 的分布, 而这就是我们谈及高斯过程所要说的。我们知道, 如果 $p(\mathbf{w})$ 是高斯的, 则每个 y 都是高斯的线性组合, 并且也是高斯的; 特殊地, 我们对 N 个输入点 $\mathbf{x}^i (i=1, \dots, N)$ 上计算的 y 值的联合分布感兴趣 (Mackay 1998)。

我们假定有 0 均值的高斯先验

$$p(\mathbf{w}) \sim \mathcal{N}(0, \alpha^{-1} \mathbf{I})$$

给定 $N \times d$ 的数据点 \mathbf{X} 和 $d \times 1$ 的权重向量, 我们将输出 y 写作

$$\mathbf{y} = \mathbf{X}\mathbf{w} \quad (14.24)$$

这是 N 元高斯, 满足

$$E[\mathbf{y}] = \mathbf{X}E[\mathbf{w}] = 0 \quad (14.25)$$

$$\text{Cov}(\mathbf{y}) = E[\mathbf{y}\mathbf{y}^T] = \mathbf{X}E[\mathbf{w}\mathbf{w}^T]\mathbf{X}^T = \frac{1}{\alpha}\mathbf{X}\mathbf{X}^T \equiv \mathbf{K}$$

其中 \mathbf{K} 是格拉姆 (Gram) 矩阵, 其元素是

$$K_{i,j} \equiv K(\mathbf{x}^i, \mathbf{x}^j) = \frac{(\mathbf{x}^i)^T \mathbf{x}^j}{\alpha}$$

在高斯过程的文献中, 这称作协方差函数 (covariance function), 并且其思想与核函数相同: 如果我们使用基函数 $\phi(\mathbf{x})$ 的集合, 则把原输入的点积通过核

$$K_{i,j} = \frac{\phi(\mathbf{x}^i)^T \phi(\mathbf{x}^j)}{\alpha}$$

推广为基函数的点积。

实际的观测输出 r 由加上噪声的直线 $r = y + \varepsilon$ 给出, 其中 $\varepsilon \sim \mathcal{N}(0, \beta^{-1})$ 。对于所有 N 个数据点, 我们将它记作

$$\mathbf{r} \sim \mathcal{N}_N(0, \mathbf{C}_N), \text{ 其中 } \mathbf{C}_N = \beta^{-1} \mathbf{I} + \mathbf{K} \quad (14.26)$$

为了做出预测, 我们将新数据看作第 $N+1$ 个数据点对 (\mathbf{x}', r') , 并使用所有 $N+1$ 个数据点表示联合分布。我们有

$$\mathbf{r}_{N+1} \sim \mathcal{N}_N(0, \mathbf{C}_{N+1}) \quad (14.27)$$

其中

$$\mathbf{C}_{N+1} = \begin{bmatrix} \mathbf{C}_N & \mathbf{k} \\ \mathbf{k}^T & c \end{bmatrix}$$

其中 \mathbf{k} 是 $K(\mathbf{x}', \mathbf{x}^i) (i=1, \dots, N)$ 的 $N \times 1$ 维向量, 而 $c = K(\mathbf{x}', \mathbf{x}') + \beta^{-1}$ 。于是, 为了做出预测, 我们计算 $p(\mathbf{r}' | \mathbf{x}', \mathbf{X}, \mathbf{r})$, 它是高斯的, 满足

$$E[\mathbf{r}' | \mathbf{x}'] = \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{r}$$

$$\text{Var}(\mathbf{r}' | \mathbf{x}') = c - \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{k}$$

图 14-6 给出了一个例子, 其中使用了线性、二次和高斯核。前两个定义为它们对应的基函数的点积; 高斯核直接定义为

$$K_G(\mathbf{x}^i, \mathbf{x}^j) = \exp \left[- \frac{\|\mathbf{x}^i - \mathbf{x}^j\|^2}{s^2} \right]$$

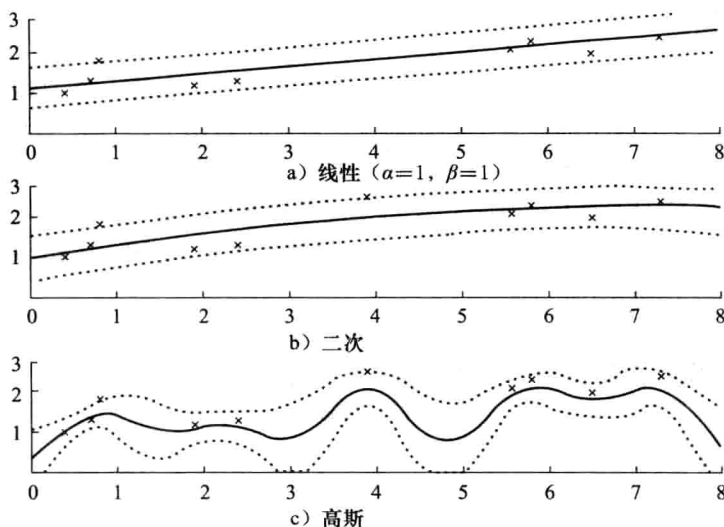


图 14-6 具有一个标准误差条线的高斯过程回归: a) 线性核, b) 二次核, c) 具有展宽 $s^2=0.5$ 的高斯核

均值是点估计(如果我们不在整个分布上积分的话), 也可以写成核效果的加权和

$$E[\mathbf{r}' | \mathbf{x}'] = \sum_t a^t K(\mathbf{x}', \mathbf{x}^t) \quad (14.28)$$

其中, a^t 是 $\mathbf{C}_N^{-1} \mathbf{r}$ 的第 t 个分量。我们还可以将它表示成训练数据点输出的加权和, 其中权

重由如下核函数给出

$$E[r' | \mathbf{x}'] = \sum_t r^t w^t \quad (14.29)$$

其中, w^t 是 $\mathbf{k}^T \mathbf{C}_N^{-1}$ 的第 t 个分量。

注意, 我们也可以计算一个点上的预测方差, 以了解其中的不确定性, 并且这依赖于影响其预测的实例。在高斯核的情况下, 只有局部区域内的那些实例是有影响的, 并且附近数据点很少的地方的预测方差高(参见图 14-7)。

正如我们在第 13 章介绍核机器时讨论的, 还可以根据应用定义和使用核函数。直接使用核函数而不必计算或存储基函数提供了很大的灵活性。通常, 给定训练集, 我们先计算参数(例如, 使用式(14.12)), 然后使用参数用式(14.13)做预测, 而不再需要训练集。这是有意义的, 因为参数的维度通常为 $\mathcal{O}(d)$, 一般比训练集的规模 N 小得多。

然而, 当我们使用基函数时, 显式计算参数可能不再如此, 因为基函数的维度可能很高, 甚至无限。在这种情况下, 正如我们这里所做的那样, 使用核函数考虑训练实例的影响, 用对偶表示可能更经济。这种思想也用于非参数光滑(参见第 8 章)和核机器(参见第 13 章)。

这里要求 \mathbf{C}_N 是可逆的, 因而是正定的。为此, \mathbf{K} 应当是半正定的, 使得将 $\beta^{-1} > 0$ 加到对角线上之后我们得到正定性。我们还看到, 最昂贵的操作是计算 $N \times N$ 矩阵的逆。幸运的是, 它只需要(在训练时)计算和存储一次。然而, 对于很大的 N , 可能还是需要近似计算。

当我们使用它对两类问题分类时, 输出要经过 S 形函数过滤, 即 $y = \text{sigmoid}(\mathbf{w}^T \mathbf{x})$, 并且 y 的分布不再是高斯的。求导类似, 不同之处也是条件概率 $p(r_{N+1} | \mathbf{x}_{N+1}, \mathbf{X}, \mathbf{r})$ 不再是高斯的, 并且需要近似计算, 例如, 使用拉普拉斯近似(Bishop 2006; Rasmussen 和 Williams 2006)。

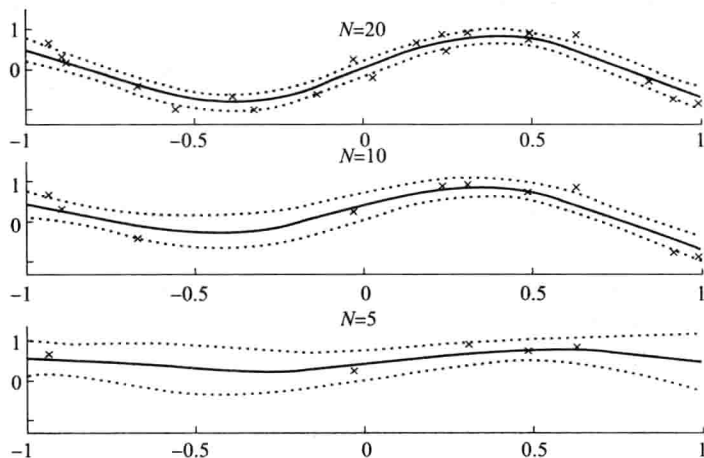


图 14-7 使用 $s^2 = 0.5$ 和不同数量的训练数据的高斯核的高斯过程回归。我们看到, 在数据不多的地方预测的方差较大

14.5 注释

随着计算能力的提高, 允许我们从后验概率抽样或近似它, 贝叶斯方法近来日趋流行。真相有多种表现形式。许多情况都偏爱简洁性, 如贝叶斯方法、正则化、最小描述长度和光滑。这是统计推断的核心思想, 因此也是机器学习的核心思想。

另一方面, 先验的主观性也是令人不安的, 并且也是反对贝叶斯方法的理由; 例如,

Gelman 2008。平坦的先验有何用，如果我们有一个尖峰的先验，为何还要收集数据？共轭先验是真实的还是仅仅为了方便？

正如支持向量机，在高斯过程中也同样如此：存在许多方法，借助它们，我们可以将新核构造其他核的函数（例如，加权和），并且这些权重或核参数（例如，展宽）可以被 2 型最大化似然过程（type 2 maximum likelihood procedure）优化，之所以称为 2 型最大化似然过程是因为我们现在不是优化参数（上面的 a' 或 w' ），而是优化第二层的超参数（Bishop 2006, Rasmussen 和 Williams 2006）。

14.6 习题

1. 对于图 14-3 的情况，观察后验如何随 N 、 σ 和 σ_0^2 变化。
2. 设 x 表示从 n 个随机样本中接收的垃圾邮件数。假定垃圾邮件所占的比例 q 的先验是 $[0, 1]$ 中的均匀分布。找出 $p(q|x)$ 的后验分布。
3. 如上，但假定 $p(q) \sim \mathcal{N}(\mu_0, \sigma_0^2)$ 。同时假定 n 很大，使得我们可以使用中心极限定理，并用高斯近似二元分布。推导 $p(q|x)$ 。
4. 在使用最大似然估计时， $\text{Var}(r')$ 是什么？将它与式 (14.17) 比较。
5. 在图 14-6 中，当 s^2 改变时拟合如何变化？
6. 提出一种过滤算法，以选择高斯过程训练集的子集。
7. 在主动学习（active learning）中，学习程序在学习期间能够逐一产生 x ，并请求监督者提供对应的 r 值，而不是被动地接受一个给定的训练集。如何用高斯过程实现主动学习？（提示：何处具有最大的不确定性？）
8. 假定我们有来自两种不同表示的输入。在这种情况下，如何使用本章讨论的方法？

14.7 参考文献

- Bishop, C. M. 2006. *Pattern Recognition and Machine Learning*. New York: Springer.
- Figueiredo, M. A. T. 2003. "Adaptive Sparseness for Supervised Learning." *IEEE Transactions on Pattern Analysis and Machine Intelligence* 25: 1150-1159.
- Gelman, A. 2008. "Objections to Bayesian statistics." *Bayesian Statistics* 3: 445-450.
- MacKay, D. J. C. 1998. "Introduction to Gaussian Processes." In *Neural Networks and Machine Learning*, ed. C. M. Bishop, 133-166. Berlin: Springer.
- MacKay, D. J. C. 2003. *Information Theory, Inference, and Learning Algorithms*. Cambridge, UK: Cambridge University Press.
- Rasmussen, C. E., and C. K. I. Williams. 2006. *Gaussian Processes for Machine Learning*. Cambridge, MA: MIT Press.
- Tibshirani, R. 1996. "Regression Shrinkage and Selection via the Lasso." *Journal of the Royal Statistical Society B* 58: 267-288.

隐马尔可夫模型

我们放松样本实例相互独立的假设，并引入马尔可夫模型，将输入序列建模为由一个参数化随机过程所产生的序列。我们讨论如何完成这一建模以及从样本序列学习模型参数的算法。

15.1 引言

迄今为止，我们一直假设样本中的实例是独立同分布的。这样做的好处是样本的似然可简化为各个实例的似然之积。然而，这一假设在相继实例相互依赖的应用中并不成立。例如，在一个单词中相继的字母是相互依赖的；在英文中，“h”非常可能跟随在“t”而非“x”之后。这类存在观测序列（例如，单词中的字母，DNA 序列中的基对）的过程并不能用简单的概率分布进行建模。一个类似的例子是语言识别，其中语音片段由称为音素的语音基元组成；只有某些音素序列是合法的，即该语言的单词。在更高的层次上，以某种序列书写或读出单词，形成由该语言的语法和语义规则定义的语句。

一个序列可用一个参数化的随机过程（parametric random process）来刻画。本章，我们讨论如何完成这种建模以及如何从样本序列的训练集中学习模型的参数。

362
363

15.2 离散马尔可夫过程

考虑一个系统，其在任意时刻处于 N 个离散状态中的一个： S_1, S_2, \dots, S_N 。在时刻 t 的状态记作 q_t ， $t = 1, 2, \dots$ 。例如， $q_t = S_i$ 表示在时刻 t 系统处于状态 S_i 。尽管我们用“时刻”好像这应该是一个时间序列，但是这种方法对任意序列，无论是时间、空间、DNA 串上位置等，都是有效的。

系统在有规律的、间隔的离散时刻，根据以前的状态值，以给定的概率转移到一个状态：

$$P(q_{t+1} = S_j | q_t = S_i, q_{t-1} = S_k, \dots)$$

对于一阶马尔可夫模型（Markov model）的特例，系统在时刻 $t+1$ 的状态仅仅依赖于在时刻 t 的状态，而与之之前的状态无关：

$$P(q_{t+1} = S_j | q_t = S_i, q_{t-1} = S_k, \dots) = P(q_{t+1} = S_j | q_t = S_i) \quad (15.1)$$

这相当于说，给定当前的状态，未来的系统状态独立于过去的状态。这恰是谚语“今天是你余生的第一天”的数学表达版本。

我们进一步简化模型，假定转移概率（transition probability）是独立于时间的：

$$a_{ij} \equiv P(q_{t+1} = S_j | q_t = S_i) \quad (15.2)$$

满足

$$a_{ij} \geq 0 \text{ 并且 } \sum_{j=1}^N a_{ij} = 1 \quad (15.3)$$

因此, 从状态 S_i 到状态 S_j 的状态转移总是具有相同的概率, 无论这个转移在观测序列中的何时何地发生。 $\mathbf{A} = [a_{ij}]$ 是一个 $N \times N$ 的矩阵, 其每行之和均为 1。

这可看作是一个随机自动机 (stochastic automation) (见图 15-1)。从每个状态 S_i , 系统以概率 a_{ij} 转移到状态 S_j , 并且这一概率在任何时刻 t 均相同。唯一的特例是第一个状态。我们定义初始概率 (initial probability) π_i , 表示序列的第一个状态是 S_i 的概率:

$$\pi_i \equiv P(q_1 = S_i) \quad (15.4)$$

满足

$$\sum_{i=1}^N \pi_i = 1 \quad (15.5)$$

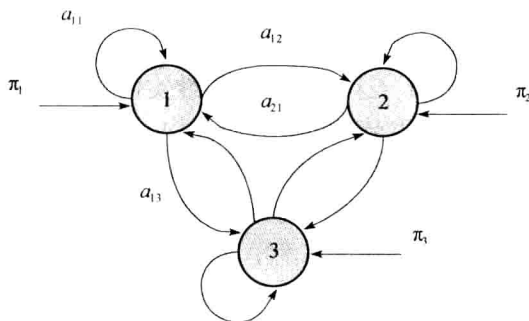


图 15-1 具有三个状态的马尔可夫模型作为随机自动机的实例。 π_i 是系统始于状态 S_i 的概率, a_{ij} 是系统从状态 S_i 转移到状态 S_j 的概率

$\Pi = [\pi_i]$ 是一个具有 N 个元素的向量, 元素和为 1。

在一个可观测马尔可夫模型 (observable Markov model) 中, 状态是可观测的。在任意时刻 t , 我们知道 q_t , 并且随着系统从一个状态转移到另一个状态, 我们得到一个观测序列, 即为状态序列。该过程的输出是一个状态集, 在每个时间点上每个状态对应一个物理可观测事件。

我们有一个观测序列 O , 它是状态序列 $O = Q = \{q_1 q_2 \cdots q_T\}$, 其概率为

$$P(O = Q | \mathbf{A}, \Pi) = P(q_1) \prod_{t=2}^T P(q_t | q_{t-1}) = \pi_{q_1} a_{q_1 q_2} \cdots a_{q_{T-1} q_T} \quad (15.6)$$

π_{q_1} 是首状态为 q_1 的概率, $a_{q_1 q_2}$ 是从 q_1 到 q_2 的概率, 以此类推。我们将这些概率相乘, 得到产生整个序列的概率。

我们现在看一个具体例子 (Rabiner 和 Juang 1986): 假定我们有 N 个容器, 每个容器中仅有一种颜色的球。例如, 有一个装红色球的容器, 一个装蓝色球的容器……。某人一个接一个地从各个容器中取球, 并将它们的颜色展示给我们。以 q_t 表示在时刻 t 所取球的颜色。我们假定有三个状态:

S_1 : 红; S_2 : 蓝; S_3 : 绿

并有初始概率:

$$\Pi = [0.5, 0.2, 0.3]^T$$

a_{ij} 是从容器 i 中取一个颜色 i 的球之后, 从容器 j 中取 (一个颜色 j 的) 球的概率。例如,

转移矩阵为:

$$\mathbf{A} = \begin{bmatrix} 0.4 & 0.3 & 0.3 \\ 0.2 & 0.6 & 0.2 \\ 0.1 & 0.1 & 0.8 \end{bmatrix}$$

给定 Π 和 \mathbf{A} , 很容易产生 K 个长度为 T 的随机序列。我们来看如何计算一个序列的概率: 假定前四个球是“红, 红, 绿, 绿”。这对应观测序列 $O = \{S_1, S_1, S_3, S_3\}$ 。其概率为

$$\begin{aligned} P(O | \mathbf{A}, \Pi) &= P(S_1) \cdot P(S_1 | S_1) \cdot P(S_3 | S_1) \cdot P(S_3 | S_3) \\ &= \pi_1 \cdot a_{11} \cdot a_{13} \cdot a_{33} \\ &= 0.5 \times 0.4 \times 0.3 \times 0.8 = 0.048 \end{aligned} \quad (15.7)$$

现在, 我们来看如何可以学习得到参数 Π 和 \mathbf{A} : 给定 K 个长度为 T 的序列, q_t^k 表示序列 k 在 t 时刻的状态, 初始概率的估计是以 S_i 起始的序列个数除以序列总数:

$$\hat{\pi}_i = \frac{\#\{\text{以 } S_i \text{ 开始的序列}\}}{\#\{\text{序列}\}} = \frac{\sum_k 1(q_1^k = S_i)}{K} \quad (15.8)$$

其中 $1(b)$ 为 1 如果 b 取真值, 否则 $1(b)$ 为 0。

至于转移概率, 对 a_{ij} 的估计为从 S_i 转移到 S_j 的个数除以所有序列中从 S_i 转移的总数:

$$\hat{a}_{ij} = \frac{\#\{\text{从 } S_i \text{ 到 } S_j \text{ 的转移}\}}{\#\{\text{从 } S_i \text{ 的转移}\}} = \frac{\sum_k \sum_{t=1}^{T-1} 1(q_t^k = S_i \text{ and } q_{t+1}^k = S_j)}{\sum_k \sum_{t=1}^{T-1} 1(q_t^k = S_i)} \quad (15.9)$$

\hat{a}_{12} 是一个蓝色球跟在一个红色球之后的次数除以所有序列中红色球的总数。

366

15.3 隐马尔可夫模型

在隐马尔可夫模型 (Hidden Markov Model, HMM) 中, 系统状态是不可观测的, 但是到达一个状态时, 可以记录一个观测, 这个观测是该状态的一个概率函数。我们假定每个状态的一个离散观测取自集合 $\{v_1, v_2, \dots, v_M\}$:

$$b_j(m) \equiv P(O_t = v_m | q_t = S_j) \quad (15.10)$$

$b_j(m)$ 是系统处于状态 S_j 时, 我们观测到 v_m ($m = 1, \dots, M$) 的观测概率 (observation probability) 或发射概率 (emission probability)。我们再次假定齐次模型, 其中发射概率不依赖于时间 t 。观测到的一系列 v_m 便形成了观测序列 O 。状态序列 Q 是不可观测的, 这正是称之为“隐”模型的缘由, 但是状态序列可以通过观测序列推断。注意, 通常许多不同的状态序列 Q 可以产生相同的观测序列 O , 但是以不同的概率产生; 正如给定服从正态分布的一个独立同分布 (iid) 的样本, 有无限多对可能的 (μ, σ) 值, 我们感兴趣的是能以最大似然产生这个样本的那对 (μ, σ) 。

还需要注意的是, 在隐马尔可夫模型中, 随机性源自两个方面: 除了从一个状态转移到另一状态是随机的之外, 系统在一个状态中产生的观测也是随机的。

再次回到我们的例子: 隐马尔可夫模型对应的容器-球实例中, 每个容器包含不同颜色的球。以 $b_j(m)$ 表示从容器 j 取出一个 m 颜色球的概率。我们再次得到一个球颜色的观测序

列, 但并不知道抽取球的容器序列。因此好像容器置于一个布帘之后, 一个人随机地从一个容器中取一个球, 而展示给我们的仅仅是球而不展示抽取球的容器。球展示后被放回容器以保持发射概率不变。球的颜色数可能不同于容器数。例如, 我们假定有三个容器, 而观测序列为:

$$O = \{\text{红, 红, 绿, 蓝, 黄}\}$$

在前面的情况下, 知道观测(球的颜色), 我们可以确切知道系统状态(容器), 因为对不同的颜色的球有不同的容器, 而且每个容器只含有一种颜色的球。可观测马尔可夫模型是隐马尔可夫模型的一个特例, 其中 $M=N$, 并且如果 $j=m$, $b_j(m)$ 为 1, 否则 $b_j(m)$ 为 0。但是在隐马尔可夫模型中, 一个球可能取自任意容器。在这种情况下, 对于相同的观测序列 O , 可能存在多个可能的状态序列 Q 产生 O (见图 15-2)。

367

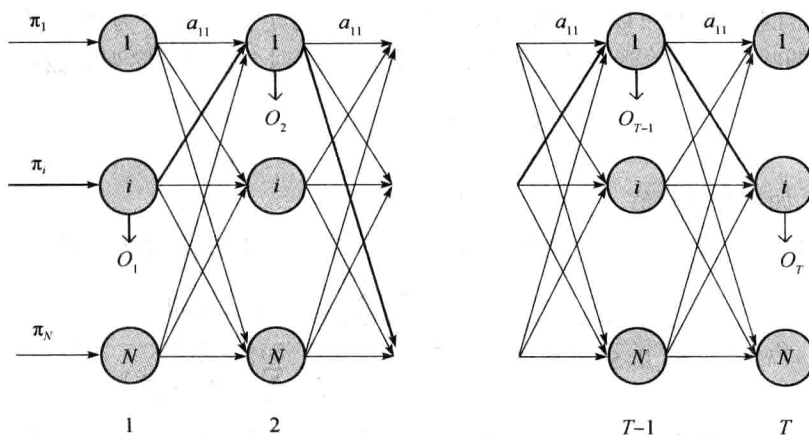


图 15-2 一个 HMM 按时间展开为格, 展示了所有可能的轨道。以粗线所示的一条路径是产生观测序列的真正(未知)状态轨迹

对上述进行总结和形式化, 一个 HMM 具有以下元素:

1) N : 模型状态个数

$$S = \{S_1, S_2, \dots, S_N\}$$

2) M : 以字母序排列的不同观测符号的个数

$$V = \{v_1, v_2, \dots, v_M\}$$

3) 状态转移概率:

$$\mathbf{A} = [a_{ij}], \quad \text{其中 } a_{ij} \equiv P(q_{t+1} = S_j | q_t = S_i)$$

4) 观测概率:

$$\mathbf{B} = [b_j(m)], \quad \text{其中 } b_j(m) \equiv P(O_t = v_m | q_t = S_j)$$

368

5) 初始状态概率:

$$\Pi = [\pi_i], \quad \text{其中 } \pi_i \equiv P(q_1 = S_i)$$

N 和 M 隐含地定义在其他参数中, 因此 $\lambda = (\mathbf{A}, \mathbf{B}, \Pi)$ 被取作一个 HMM 的参数集。

给定 λ , 模型可用于产生任意长度的任意个数观测序列, 但是我们通常感兴趣的是另一方向, 即通过一个由若干观测序列组成的训练集来估计模型的参数。

15.4 HMM 的三个基本问题

给定一定数量的观测序列，我们对下面三个问题感兴趣：

1) 给定一个模型 λ ，我们希望估计任意给定观测序列 $O = \{O_1 O_2 \cdots O_T\}$ 的概率，即 $P(O | \lambda)$ 。

2) 给定一个模型 λ 及一个观测序列 O ，我们希望找出状态序列 $Q = \{q_1 q_2 \cdots q_T\}$ ，其具有产生 O 的最大概率，即我们要找到最大化 $P(Q | O, \lambda)$ 的 Q^* 。

3) 给定观测序列组成的训练集 $X = \{O^k\}_k$ ，我们希望学习得到产生 X 的概率最大化的模型，即我们要找到最大化 $P(X | \lambda)$ 的 λ^* 。

让我们逐一来看这些问题的解决办法，每个办法又用于解决下一个问题，直到我们得以计算出 λ ，或者是从数据中学到一个模型。

15.5 估值问题

给定观测序列 $O = \{O_1 O_2 \cdots O_T\}$ 和状态序列 $Q = \{q_1 q_2 \cdots q_T\}$ ，给定状态序列 Q 观测到 O 的概率为

$$P(O | Q, \lambda) = \prod_{t=1}^T P(O_t | q_t, \lambda) = b_{q_1}(O_1) \cdot b_{q_2}(O_2) \cdots b_{q_T}(O_T) \quad (15.11) \quad \boxed{369}$$

我们无法直接计算上式，因为状态序列(的值)是未知的。状态序列 Q 的概率是

$$P(Q | \lambda) = P(q_1) \prod_{t=2}^T P(q_t | q_{t-1}) = \pi_{q_1} a_{q_1 q_2} \cdots a_{q_{T-1} q_T} \quad (15.12)$$

因而联合概率为

$$\begin{aligned} P(O, Q | \lambda) &= P(q_1) \prod_{t=2}^T P(q_t | q_{t-1}) \prod_{t=1}^T P(O_t | q_t) \\ &= \pi_{q_1} b_{q_1}(O_1) a_{q_1 q_2} b_{q_2}(O_2) \cdots a_{q_{T-1} q_T} b_{q_T}(O_T) \end{aligned} \quad (15.13)$$

我们可以通过边缘化这一联合概率，即通过在所有可能的 Q 上求和，来计算 $P(O | \lambda)$ ：

$$P(O | \lambda) = \sum_{\text{所有可能的 } Q} P(O, Q | \lambda)$$

但是，该方法是不现实的，因为如果假定所有的概率都是非零的，则有 N^T 个可能的 Q 。幸运的是，存在计算 $P(O | \lambda)$ 的有效方法，称之为正反向过程(forward-backward procedure)，见图 15-3。其基于的思想是将观测序列分为两个部分：第一部分始于时刻 1 到时刻 t ，第二部分始于时刻 $t+1$ 直到 T 。

给定模型 λ ，我们定义正向变量(forward variable) $\alpha_t(i)$ 为到时刻 t 观测到部分序列 $\{O_1 \cdots O_t\}$ 并且在时刻 t 的状态为 S_i 时的概率：

$$\alpha_t(i) \equiv P(O_1 \cdots O_t, q_t = S_i | \lambda) \quad (15.14)$$

这种方法的优点在于可通过结果累积而递归地计算上式：

■ 初始化：

$$\begin{aligned} \alpha_1(i) &\equiv P(O_1, q_1 = S_i | \lambda) \\ &= P(O_1 | q_1 = S_i, \lambda) P(q_1 = S_i | \lambda) \\ &= \pi_i b_i(O_1) \end{aligned} \quad (15.15)$$

■ 递归(见图 15-3(a)):

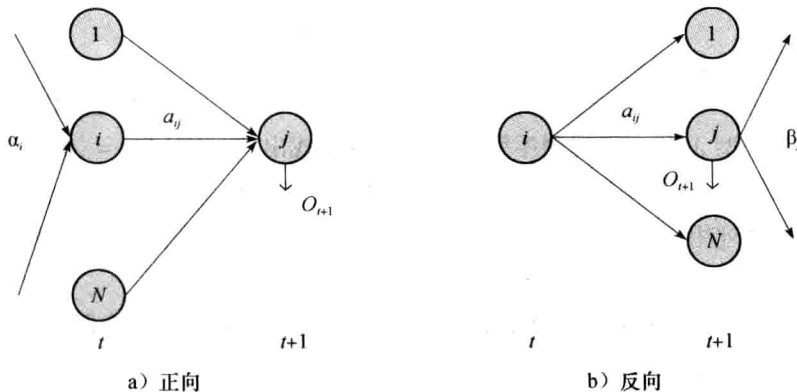


图 15-3 正反向过程 a) $\alpha_t(j)$ 的计算; b) $\beta_t(i)$ 的计算

$$\begin{aligned}
 \alpha_{t+1}(j) &\equiv P(O_1 \cdots O_{t+1}, q_{t+1} = S_j | \lambda) \\
 &= P(O_1 \cdots O_{t+1} | q_{t+1} = S_j, \lambda) P(q_{t+1} = S_j | \lambda) \\
 &= P(O_1 \cdots O_t | q_{t+1} = S_j, \lambda) P(O_{t+1} | q_{t+1} = S_j, \lambda) P(q_{t+1} = S_j | \lambda) \\
 &= P(O_1 \cdots O_t, q_{t+1} = S_j | \lambda) P(O_{t+1} | q_{t+1} = S_j, \lambda) \\
 &= P(O_{t+1} | q_{t+1} = S_j, \lambda) \sum_i P(O_1 \cdots O_t, q_t = S_i, q_{t+1} = S_j, \lambda) \\
 &= P(O_{t+1} | q_{t+1} = S_j, \lambda) \sum_i P(O_1 \cdots O_t, q_{t+1} = S_j | q_t = S_i, \lambda) P(q_t = S_i | \lambda) \\
 &= P(O_{t+1} | q_{t+1} = S_j, \lambda) \sum_i P(O_1 \cdots O_t | q_t = S_i, \lambda) P(q_{t+1} = S_j | q_t = S_i, \lambda) P(q_t = S_i | \lambda) \\
 &= P(O_{t+1} | q_{t+1} = S_j, \lambda) \sum_i P(O_1 \cdots O_t, q_t = S_i | \lambda) P(q_{t+1} = S_j | q_t = S_i, \lambda) \\
 &= \left[\sum_{i=1}^N \alpha_t(i) a_{ij} \right] b_j(O_{t+1}) \tag{15.16}
 \end{aligned}$$

$\alpha_t(i)$ 表示了产生前 t 个观测并且止于状态 S_i 的概率。通过将其乘以概率 a_{ij} 得到转移到状态 S_j 的概率, 但是因为有 N 个可能的前驱状态, 我们需要对所有可能的前驱状态 S_i 求和。 $b_j(O_{t+1})$ 则是产生第 $(t+1)$ 个观测且在第 $t+1$ 时刻处于状态 S_j 的概率。

当我们计算正向变量时, 容易计算观测序列的概率:

$$\begin{aligned}
 P(O | \lambda) &= \sum_{i=1}^N P(O, q_T = S_i | \lambda) \\
 &= \sum_{i=1}^N \alpha_T(i) \tag{15.17}
 \end{aligned}$$

$\alpha_T(i)$ 是产生整个观测序列并终止于状态 S_i 的概率。我们需要对所有可能的终止状态进行求和。

计算 $\alpha_t(i)$ 的复杂度为 $\mathcal{O}(N^2 T)$, 并且这在合理的时间内解决我们第一个估值问题。虽然现在我们并不需要, 但是我们类似地定义反向变量 (backward variable) $\beta_t(i)$, 作为在时刻 t 处于状态 S_i 并且观测到部分序列 $O_{t+1} \cdots O_T$ 的概率:

$$\beta_i(i) \equiv P(O_{t+1} \cdots O_T | q_t = S_i, \lambda) \quad (15.18)$$

同样可以按如下进行递归计算，不同的是这次是逆序进行：

- 初始化(任意地设为1)：

$$\beta_T(i) = 1$$

- 递归计算(见图 15-3(b))：

$$\begin{aligned} \beta_i(i) &\equiv P(O_{t+1} \cdots O_T | q_t = S_i, \lambda) \\ &= \sum_j P(O_{t+1} \cdots O_T, q_{t+1} = S_j | q_t = S_i, \lambda) \\ &= \sum_j P(O_{t+1} \cdots O_T | q_{t+1} = S_j, q_t = S_i, \lambda) P(q_{t+1} = S_j | q_t = S_i, \lambda) \\ &= \sum_j P(O_{t+1} | q_{t+1} = S_j, q_t = S_i, \lambda) \\ &\quad P(O_{t+2} \cdots O_T | q_{t+1} = S_j, q_t = S_i, \lambda) P(q_{t+1} = S_j | q_t = S_i, \lambda) \\ &= \sum_j P(O_{t+1} | q_{t+1} = S_j, \lambda) \\ &\quad P(O_{t+2} \cdots O_T | q_{t+1} = S_j, \lambda) P(q_{t+1} = S_j | q_t = S_i, \lambda) \\ &= \sum_{j=1}^N a_{ij} b_j(O_{t+1}) \beta_{t+1}(j) \end{aligned} \quad (15.19)$$

372

当处于状态 S_i 时，有 N 种可能的下一状态 S_j ，每个的概率为 a_{ij} 。在该状态上，我们产生第 $(t+1)$ 个观测，而 $\beta_{t+1}(j)$ 表示了时刻 $t+1$ 之后的所有观测的概率，以此类推。

对于实现，需要引起注意的是： α_t 和 β_t 都是通过多个小概率相乘计算的，而当序列很长时有下溢的危险。为了避免下溢，我们在每一步通过将 $\alpha_t(i)$ 乘以

$$c_t = \frac{1}{\sum_j \alpha_t(j)}$$

对其进行规范化。同样也将 $\beta_t(i)$ 乘以相同的 c_t ，以对其进行规范化($\beta_t(i)$ 之和不为 1)。规范化之后不能使用式(15.17)；替代地，我们有(Rabiner 1989)

$$P(O | \lambda) = \frac{1}{\prod_t c_t} \quad \text{或} \quad \log P(O | \lambda) = - \sum_t \log c_t \quad (15.20)$$

15.6 寻找状态序列

我们现在考虑第二个问题，即给定模型 λ ，寻找以最高的概率产生观测序列 $O = \{O_1, O_2, \cdots, O_T\}$ 的状态序列 $Q = \{q_1, q_2, \cdots, q_T\}$ 。

定义 $\gamma_t(i)$ 为给定 O 和 λ ，在时刻 t 处于状态 S_i 的概率，它可以计算如下：

$$\begin{aligned} \gamma_t(i) &\equiv P(q_t = S_i | O, \lambda) \\ &= \frac{P(O | q_t = S_i, \lambda) P(q_t = S_i | \lambda)}{P(O | \lambda)} \\ &= \frac{P(O_1 \cdots O_t | q_t = S_i, \lambda) P(O_{t+1} \cdots O_T | q_t = S_i, \lambda) P(q_t = S_i | \lambda)}{\sum_{j=1}^N P(O, q_t = S_j | \lambda)} \end{aligned} \quad (15.21)$$

$$\begin{aligned}
&= \frac{P(O_1 \cdots O_t, q_t = S_i | \lambda) P(O_{t+1} \cdots O_T | q_t = S_i, \lambda)}{\sum_{j=1}^N P(O | q_t = S_j, \lambda) P(q_t = S_j | \lambda)} \\
&= \frac{\alpha_t(i) \beta_t(i)}{\sum_{j=1}^N \alpha_t(j) \beta_t(j)} \quad (15.22)
\end{aligned}$$

这里我们看到 $\alpha_t(i)$ 和 $\beta_t(i)$ 是如何很好的在它们之间分割序列：正向变量 $\alpha_t(i)$ 解释了到时刻 t 为止的序列的前一部分并终止于状态 S_i ，而反向变量 $\beta_t(i)$ 从那里开始并解释了直到时刻 T 的后一部分。

分子 $\alpha_t(i) \beta_t(i)$ 表示了时刻 t 、系统处于状态 S_i 的整个序列的概率。我们需要将其除以所有在时刻 t 可能转移到的中间状态对其进行正规化，并保证 $\sum_i \gamma_t(i) = 1$ 。

为找到(目标)状态序列，可以在每一步 t 选择具有最高概率的状态：

$$q_t^* = \arg \max_i \gamma_t(i) \quad (15.23)$$

但是这有可能在时刻 t 和时刻 $t+1$ 选择 S_i 和 S_j 作为最合适状态，即使这时有 $a_{ij}=0$ 。为了找到单个最好的状态序列(路径)，我们使用基于动态规划的 Viterbi 算法(Viterbi algorithm)，将这样的转移概率考虑在内。

给定状态序列 $Q = q_1 q_2 \cdots q_T$ 和观测序列 $O = O_1 \cdots O_T$ ，定义 $\delta_t(i)$ 为在时刻 t 代表前 t 个观测并止于状态 S_i 具有最高概率路径的概率：

$$\delta_t(i) \equiv \max_{q_1 q_2 \cdots q_{t-1}} p(q_1 q_2 \cdots q_{t-1}, q_t = S_i, O_1 \cdots O_t | \lambda) \quad (15.24)$$

从而我们可以递归地计算 $\delta_{t+1}(i)$ ，而最优路径可以从时刻 T 开始，在每个时刻选择最合适的状态来反向读取。算法如下：

1) 初始化：

$$\delta_1(i) = \pi_i b_i(O_1)$$

$$\psi_1(i) = 0$$

2) 递归：

$$\delta_t(j) = \max_i \delta_{t-1}(i) a_{ij} \cdot b_j(O_t)$$

$$\psi_t(j) = \arg \max_i \delta_{t-1}(i) a_{ij}$$

3) 终止：

$$p^* = \max_i \delta_T(i)$$

$$q_T^* = \arg \max_i \delta_T(i)$$

4) 路径(状态序列)回溯：

$$q_t^* = \psi_{t+1}(q_{t+1}^*), t = T-1, T-2, \dots, 1$$

使用图 15-2 的网格结构， $\psi_t(j)$ 跟踪了在时刻 $t-1$ 最大化 $\delta_t(j)$ 的状态，即最佳的前驱状态。Viterbi 算法与正向阶段具有相同的复杂度，其中我们在每一步用取最大值替代求和。

15.7 学习模型参数

我们现在继续第三个问题，从数据中学习 HMM。使用最大似然方法，我们要计算最大化训练序列样本 $\chi = \{O^k\}_{k=1}^K$ 的似然的 λ^* ，即计算最大化 $P(\chi | \lambda)$ 的 λ^* 。我们从定义便于稍后讨论的新变量开始。

定义 $\xi_t(i, j)$ 为给定全部观测 O 和 λ ，在时刻 t 处于状态 S_i 和在时刻 $t+1$ 处于状态 S_j 的概率：

$$\xi_t(i, j) = P(q_t = S_i, q_{t+1} = S_j | O, \lambda) \quad (15.25)$$

上式可计算为(参见图 15-4)：

$$\begin{aligned} \xi_t(i, j) &= P(q_t = S_i, q_{t+1} = S_j | O, \lambda) \\ &= \frac{P(O | q_t = S_i, q_{t+1} = S_j, \lambda) P(q_t = S_i, q_{t+1} = S_j | \lambda)}{P(O | \lambda)} \\ &= \frac{P(O | q_t = S_i, q_{t+1} = S_j, \lambda) P(q_{t+1} = S_j | q_t = S_i, \lambda) P(q_t = S_i | \lambda)}{P(O | \lambda)} \\ &= \left(\frac{1}{P(O | \lambda)} \right) P(O_1 \cdots O_t | q_t = S_i, \lambda) P(O_{t+1} | q_{t+1} = S_j, \lambda) \\ &\quad P(O_{t+2} \cdots O_T | q_{t+1} = S_j, \lambda) a_{ij} P(q_t = S_i | \lambda) \\ &= \left(\frac{1}{P(O | \lambda)} \right) P(O_1 \cdots O_t, q_t = S_i | \lambda) P(O_{t+1} | q_{t+1} = S_j, \lambda) \\ &\quad P(O_{t+2} \cdots O_T | q_{t+1} = S_j, \lambda) a_{ij} \\ &= \frac{\alpha_t(i) b_j(O_{t+1}) \beta_{t+1}(j) a_{ij}}{\sum_k \sum_l P(q_t = S_k, q_{t+1} = S_l, O | \lambda)} \\ &= \frac{\alpha_t(i) a_{ij} b_j(O_{t+1}) \beta_{t+1}(j)}{\sum_k \sum_l \alpha_t(k) a_{kl} b_l(O_{t+1}) \beta_{t+1}(l)} \end{aligned} \quad (15.26)$$

$\alpha_t(i)$ 表示产生前 t 个观测并且在时刻 t 止于状态 S_i 的概率。以概率 a_{ij} 转移到态 S_j ，产生第 $t+1$ 个观测，并在 $t+1$ 时刻从 S_j 开始继续产生其余的观测序列。我们通过将 $\xi_t(i, j)$ 除以所有在时刻 t 和时刻 $t+1$ 可能处于的状态对其进行规范化。

如果需要，我们可以通过对所有可能的下一状态，在弧概率上边缘化来计算在 t 时刻系统处于状态 S_i 的概率：

$$\gamma_t(i) = \sum_{j=1}^N \xi_t(i, j) \quad (15.27)$$

需要注意的是，如果马尔可夫模型不是隐的而是可观测的，则 $\gamma_t(i)$ 和 $\xi_t(i, j)$ 二者均为 0/1。而当模型是隐马尔可夫模型时，我们通过软计数 (soft count) 这样的后验概率来估计它们。正如监督分类和非监督聚类之间的区别，其中类标号相应的为已知和未知。在使用

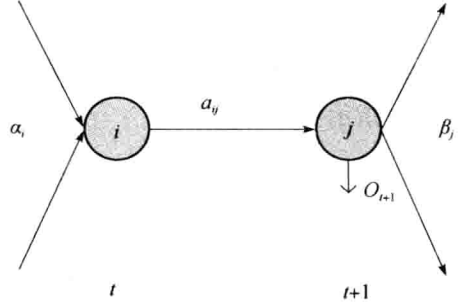


图 15-4 计算弧概率 $\xi_t(i, j)$

EM 算法的非监督聚类中(7.4 节), 类标号未知, 我们首先(在 E-步骤中)估计它们, 而后(在 M-步骤中)使用这些估计计算参数。

类似地, 这里我们使用 Baum-Welch 算法(Baum-Welch algorithm), 它是一种 EM 方法。在每次迭代中, 首先在 E-步, 我们在给定当前 $\lambda = (\mathbf{A}, \mathbf{B}, \Pi)$ 情况下计算 $\xi_t(i, j)$ 和 $\gamma_t(i)$ 的值, 然后在 M-步, 我们在给定 $\xi_t(i, j)$ 和 $\gamma_t(i)$ 的情况下再计算 λ 。这两个步骤交替进行直到收敛, 这是因为, $P(O|\lambda)$ 的值在这个过程中不会减小。

假设指示变量 z'_i 为:

$$z'_i = \begin{cases} 1 & \text{如果 } q_t = S_i \\ 0 & \text{否则} \end{cases} \quad (15.28)$$

并且

$$z'_{ij} = \begin{cases} 1 & \text{如果 } q_t = S_i \text{ 并且 } q_{t+1} = S_j \\ 0 & \text{否则} \end{cases} \quad (15.29)$$

这些值在可观测马尔可夫模型情况下为 0/1, 而在 HMM 情况下为隐随机变量。在后一种情况下, 我们在 E-步对其进行估计:

$$\begin{aligned} E[z'_i] &= \gamma_t(i) \\ E[z'_{ij}] &= \xi_t(i, j) \end{aligned} \quad (15.30)$$

在 M-步, 我们在这些估计值上计算参数。从 S_i 到 S_j 的转移的期望数为 $\sum_t \xi_t(i, j)$, 而从 S_i 转移的总数为 $\sum_t \gamma_t(i)$ 。这两个数值的比值给出了任意时刻从状态 S_i 转移到 S_j 的概率:

$$\hat{a}_{ij} = \frac{\sum_{t=1}^{T-1} \xi_t(i, j)}{\sum_{t=1}^{T-1} \gamma_t(i)} \quad (15.31)$$

注意除了将实际的计数替换为估计的软计数外, 上式和式(15.9)是一样的。

在状态 S_j 观测到 v_m 的概率为系统处于状态 S_j 时观测到 v_m 的期望次数除以系统处于状态 S_j 的总数:

$$\hat{b}_j(m) = \frac{\sum_{t=1}^T \gamma_t(j) 1(O_t = v_m)}{\sum_{t=1}^T \gamma_t(j)} \quad (15.32)$$

当有多个观测序列 $\mathcal{X} = \{O^k\}_{k=1}^K$ 时, 我们假定它们是相互独立的:

$$P(\mathcal{X}|\lambda) = \prod_{k=1}^K P(O^k|\lambda)$$

参数在全部序列的所有观测上取平均:

$$\hat{a}_{ij} = \frac{\sum_{k=1}^K \sum_{t=1}^{T_K-1} \xi_t^k(i, j)}{\sum_{k=1}^K \sum_{t=1}^{T_K-1} \gamma_t^k(i)}$$

$$\hat{b}_j(m) = \frac{\sum_{k=1}^K \sum_{t=1}^{T_K} \gamma_t^k(j) 1(O_t^k = v_m)}{\sum_{k=1}^K \sum_{t=1}^{T_K-1} \gamma_t^k(j)}$$

$$\hat{\pi}_i = \frac{\sum_{k=1}^K \gamma_1^k(i)}{K} \quad (15.33)$$

15.8 连续观测

在我们的讨论中，我们假定离散的观测服从多项分布：

$$P(O_t | q_t = S_j, \lambda) = \prod_{m=1}^M b_j(m) r_m^t \quad (15.34)$$

其中

$$r_m^t = \begin{cases} 1 & \text{如果 } O_t = v_m \\ 0 & \text{否则} \end{cases} \quad (15.35)$$

如果输入是连续的，一个方法是将其离散化，然后使用这些离散值作为观测值。通常使用向量量化(7.3节)，目的在于将连续值转换为最接近的参考向量的离散值。例如，在语音识别中，一个单词发音被分割为小的语音片段，对应音节或部分音节；预处理后，这些片段通过向量量化被离散化，进而使用 HMM 将一个单词发音建模为一个离散化片段的序列。

我们还记得用于向量量化的 k -均值是高斯混合模型的一个硬版本：

$$p(O_t | q_t = S_j, \lambda) = \sum_{l=1}^L P(G_l) p(O_t | q_t = S_j, G_l, \lambda) \quad (15.36)$$

其中

$$p(O_t | q_t = S_j, G_l, \lambda) \sim \mathcal{N}(\mu_l, \Sigma_l) \quad (15.37)$$

并且观测保持连续性。在这种高斯混合情形下，可为分量参数(以合适的正则化来保持对参数个数进行检验)和混合比例推导出 EM 方程(Rabiner 1989)。

378

现在我们看一下观测为连续标量的情形， $O_t \in \mathcal{R}$ 。最简单的方法是假定其服从正态分布：

$$p(O_t | q_t = S_j, \lambda) \sim \mathcal{N}(\mu_j, \sigma_j^2) \quad (15.38)$$

这意味在状态 S_j ，观测取自均值为 μ_j ，方差为 σ_j^2 的正态分布。这种情形下，M 步的公式为：

$$\hat{\mu}_j = \frac{\sum_t \gamma_t(j) O_t}{\sum_t \gamma_t(j)} \quad (15.39)$$

$$\hat{\sigma}_j^2 = \frac{\sum_t \gamma_t(j) (O_t - \hat{\mu}_j)^2}{\sum_t \gamma_t(j)}$$

15.9 带输入的 HMM

在某些应用中,除了观测序列 O_t 之外,我们还有一个输入序列 x_t 。我们可以将处于状态 S_j 的观测 O_t 以输入 x_t 为条件,并记作 $P(O_t | q_t = S_j, x_t)$ 。当观测为连续标量的情况下,我们将式(15.38)替换为一个广义模型

$$p(O_t | q_t = S_j, x_t, \lambda) \sim \mathcal{N}(g_j(x^t | \theta_j), \sigma_j^2) \quad (15.40)$$

其中,比如假定线性模型,我们有

$$g_j(x^t | w_j, w_{j0}) = w_j x^t + w_{j0} \quad (15.41)$$

如果观测是离散的并服从多项分布,则我们得到一个将 x_t 作为输入并产生 M 选 1 输出的分类器,否则我们可以产生后验类概率并保持观测的连续性。

类似地,状态转移概率也可以以输入为条件,即 $P(q_{t+1} = S_j | q_t = S_i, x_t)$ 。它可通过一个选择将 $t+1$ 时刻的状态作为在 t 时刻的状态和输入的函数的分类器来实现。这就是马尔可夫混合专家模型(markov mixture of expert)(Meila 和 Jordan 1996),是混合专家构架的一般化(见 12.8 节),其中门网络跟踪其在前一时步所作的决策。此类构架也称为输入-输出 HMM(input-output HMM)(Bengio 和 Frasconi 1996)并具有模型非齐次的优点;在不同的时步使用不同的观测和转移概率。在每一个状态仍然使用由 θ_j 参数化的单一模型,但是根据所看到的输入产生不同的转移或观测概率。有可能输入并非单值,而是围绕时刻 t 的一个窗口,使得输入为一个向量;这样可以处理输入和观测序列具有不同长度的应用。

即使没有其他显式的输入序列,带输入的 HMM 也可以通过关于一些以前观测的预设函数来产生一个“输入”

$$x_t = f(O_{t-\tau}, \dots, O_{t-1})$$

从而提供一个大小为 τ 的背景输入窗口。

15.10 HMM 中的模型选择

正如任意其他模型一样,需要调整 HMM 的复杂度,使其根据面对的数据的规模和性质平衡其复杂性。一种可能的方法是调整 HMM 的拓扑。在完全连接(遍历)的 HMM 中,从一个状态可转移到任意其他状态,使得 \mathbf{A} 是一个 $N \times N$ 的全矩阵。在一些应用中,仅有某些转移是允许的,而不允许的转移有 $a_{ij} = 0$ 。当下一个可能的状态较少,即 $N' < N$ 时,正反向扫描和 Viterbi 过程的复杂度为 $\mathcal{O}(NN'T)$,而并非 $\mathcal{O}(N^2T)$ 。

例如,在语音识别中使用自左向右 HMM(left-to-right HMM),其中系统状态按时间排序,随着时间的推进,状态下标增加或保持不变。这样的约束可用来对像语音一样其性质随时间变化的序列进行建模,并且当到达一个状态的时候,我们近似地知道其前的状态。有一个性质:系统绝不向具有更小下标的状态转移,即对于 $j < i$ 有 $a_{ij} = 0$ 。而在状态下标上跨度很大的状态转移也不允许,即对于 $j > i + \tau$ 有 $a_{ij} = 0$ 。图 15-5 给出了自左向右 HMM 的一个例子,其中 $\tau = 2$,状态转移矩阵如下:

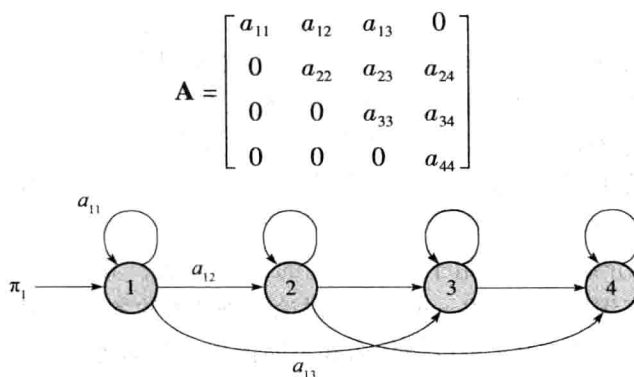


图 15-5 自左向右 HMM 的例子

决定 HMM 复杂度的另一因素是状态的个数 N 。因为状态是隐藏的，因此其个数未知并且应在训练前选定。这需要使用先验信息对其进行决定并通过交叉验证，即通过检查验证序列的似然进行微调。

当用于分类时，可使用一组 HMM，每个对属于一类的序列进行建模。例如，在口语单词识别中，每个单词的样本训练一个不同的模型 λ_i 。当对新的单词发音 O 进行分类时，所有不同的单词模型均被用来计算 $P(O | \lambda_i)$ 。然后使用贝叶斯规则得到后验概率

$$P(\lambda_i | O) = \frac{P(O | \lambda_i) P(\lambda_i)}{\sum_j P(O | \lambda_j) P(\lambda_j)} \quad (15.42)$$

其中 $P(\lambda_i)$ 是单词 i 的先验概率。该发音被指派到具有最高后验概率的单词。这是基于似然的方法，也是直接训练有判别力的 HMM，以便最大化后验概率。当存在同一单词的多个发音时，它们在该单词的 HMM 中被定义为并行路径。

在像语音这样的连续输入的情况下，困难之处在于将信号分割为较小的离散观测。通常使用取作基元部分的音素 (phone)，并通过对其组合形成更长的序列 (例如单词)。每个音素 (通过向量量化) 被并行地识别，然后用 HMM 将它们顺序组合。如果语音基元简单，则 HMM 会比较复杂，反之亦然。在连续语音识别中，单词并非一个接一个的以清晰间隔进行发音，这时可以采用多级别上的层次 HMM；一层用于组合音素以识别单词，另一层通过建立语言模型对单词组合以识别语句等。

近年来，神经网络/HMM 混合模型在语音识别领域比较流行 (Morgan 和 Bourlard 1995)。在这样的模型中，一个多层感知器 (第 11 章) 用于捕捉时间局部而不是比较复杂和非线性的基元，如音素，而 HMM 用于学习时间结构。神经网络作为预处理器并将时间窗口中的原始观测转换成比向量量化的输出更容易建模的形式。

正如我们将在第 16 章看到的，HMM 可以看作一种图模型，而 HMM 中的估计可以看作信念传播算法的一个特例。我们专门写这一章是因为这种特定模型具有广泛和成功的应用，特别是在自动语音识别领域。当我们详细讨论图模型时，将会看到如何扩展基本 HMM 架构，例如，通过有多个序列，或通过引入可以简化模型的隐藏变量来扩展基本 HMM 架构。

15.11 注释

HMM 是一项成熟的技术，并且存在基于 HMM 的商业语音识别系统在实际使用 (Rabiner

和 Juang 1993; Jelinek 1997)。在 11.12 节, 我们讨论了如何训练多层感知器用于序列识别。与延迟神经网络相比, HMM 的优点在于不用事先定义时间窗口, 并且 HMM 的训练效果优于递归神经网络。HMM 可以应用于各种序列识别任务。HMM 在生物信息领域的应用在 Baldi 和 Brunak 1998 中有所介绍, 在自然语音处理中的应用在 Manning 和 Schütz 1999 有所介绍。HMM 也用于在线手写符号识别, 它与光学识别的不同之处在于书写者在触觉感知的书写板上书写, 并且输入是一个关于笔尖在书写板上移动的 (x, y) 坐标序列, 而非静态的图片。Bengio 等(1995)介绍了一种用于在线识别的混合系统, 其中 MLP 识别单个字符, 而 HMM 将字符组合起来识别单词。Bengio 1999 讨论了 HMM 的各种应用和多种扩展, 例如有判别力的 HMM。一个关于 HMM 可以做什么和不能做什么的综述参见 Bilmes 2006。

在任意此类识别系统中, 一个关键点在于决定多少工作并行进行以及将什么工作留做串行处理。在语音识别中, 音素可通过一个并行系统来识别, 这相当于假定所有的音素声音同时发出。然后, 通过组合音素顺序识别单词。在其他系统中, 如果相同的音素有多个版本, 比如根据之前以及之后音素的不同, 音素本身可被设计为更简单的语音发声的序列。并行工作是好的, 但仅仅是在一定程度上; 我们应当在并行和顺序处理之间找到理想的平衡。为了可以一键式接通任意人的电话, 我们可能需要电话上有百万按键; 作为替代, 我们用十个按键并以序列按键来拨通电话。

我们将在第 16 章讨论图模型, 其中我们将会看到 HMM 可以看作一类特殊的图模型, 并且 HMM 上的推断和学习操作类似于贝叶斯网络中的对应操作 (Smyth, Heckerman 和 Jordan 1997)。正如我们稍后就会看到的, HMM 有多种扩展, 如因子 HMM (factorial HMM), 其中在每一时间步, 一定数量的集体生成观测的状态以及树结构 HMM (tree-structured HMM), 其中存在状态的层次关系。这种一般的形式化使得我们可以处理连续状态和离散状态, 称作线性动态系统 (linear dynamical system)。对于这类模型中的某些, 不可能做精确推断, 我们需要近似或抽样方法 (Ghahramani 2001)。

15.12 习题

1. 给定一个有三个状态 S_1, S_2, S_3 的可观测马尔可夫模型, 其初始概率为:

$$\Pi = [0.5, 0.2, 0.3]^T$$

转移概率为

$$\mathbf{A} = \begin{bmatrix} 0.4 & 0.3 & 0.3 \\ 0.2 & 0.6 & 0.2 \\ 0.1 & 0.1 & 0.8 \end{bmatrix}$$

产生 100 个有 1000 个状态的序列。

2. 使用上题中产生的数据来估计 Π 和 \mathbf{A} , 并和产生这些数据的参数进行比较。
3. 形式化一个二阶马尔可夫模型。其参数是什么? 如何对一个给定的状态序列计算其概率? 对于一个可观测模型如何学习其参数?
4. 证明任意二阶 (或更高阶) 马尔可夫模型可转化为一个一阶马尔可夫模型。
5. 一些研究者将马尔可夫模型定义为当穿越一条弧 (边) 的时候产生一个观测, 而非到达一个状态。这种模型的能力是否比我们讨论的模型更强?

6. 从一个你选择的 HMM 产生训练和确认序列。在相同的训练集上通过变化隐藏状态的个数来训练多个不同的 HMM 并计算相应的确认似然。观察确认似然如何随着状态个数的增加而变化。
7. 在式(15.38)中, 如果我们有多元观测, 那么 M-步的方程是什么?
8. 考虑容器-球的例子, 如果我们不放回地抽取, 那么会有何不同?
9. 假定在任意时刻我们都有两个来自不同字母表的观测, 例如, 每天观测两种货币的币值。如何使用 HMM 实现?
10. 如何得到增量 HMM? 增量 HMM 允许在必要时增加新的隐藏状态。

15.13 参考文献

- Baldi, P., and S. Brunak. 1998. *Bioinformatics: The Machine Learning Approach*. Cambridge, MA: MIT Press.
- Bengio, Y. 1999. "Markovian Models for Sequential Data." *Neural Computing Surveys* 2: 129-162.
- Bengio, Y., and P. Frasconi. 1996. "Input-Output HMMs for Sequence Processing." *IEEE Transactions on Neural Networks* 7: 1231-1249.
- Bengio, Y., Y. Le Cun, C. Nohl, and C. Burges. 1995. "LeRec: A NN/HMM Hybrid for On-line Handwriting Recognition." *Neural Computation* 7: 1289-1303.
- Bilmes, J. A. 2006. "What HMMs Can Do." *IEICE Transactions on Information and Systems* E89-D: 869-891.
- Ghahramani, Z. 2001. "An Introduction to Hidden Markov Models and Bayesian Networks." *International Journal of Pattern Recognition and Artificial Intelligence* 15: 9-42.
- Jelinek, F. 1997. *Statistical Methods for Speech Recognition*. Cambridge, MA: MIT Press.
- Manning, C. D., and H. Schütze. 1999. *Foundations of Statistical Natural Language Processing*. Cambridge, MA: MIT Press.
- Meila, M., and M. I. Jordan. 1996. "Learning Fine Motion by Markov Mixtures of Experts." In *Advances in Neural Information Processing Systems* 8, ed. D. S. Touretzky, M. C. Mozer, and M. E. Hasselmo, 1003-1009. Cambridge, MA: MIT Press.
- Morgan, N., and H. Bourlard. 1995. "Continuous Speech Recognition: An Introduction to the Hybrid HMM/Connectionist Approach." *IEEE Signal Processing Magazine* 12: 25-42.
- Smyth, P., D. Heckerman, and M. I. Jordan. 1997. "Probabilistic Independence Networks for Hidden Markov Probability Models." *Neural Computation* 9: 227-269.
- Rabiner, L. R. 1989. "A Tutorial on Hidden Markov Models and Selected Applications in Speech Recognition." *Proceedings of the IEEE* 77: 257-286.
- Rabiner, L. R., and B. H. Juang. 1986. "An Introduction to Hidden Markov Models." *IEEE Acoustics, Speech, and Signal Processing Magazine* 3: 4-16.
- Rabiner, L. R., and B. H. Juang. 1993. *Fundamentals of Speech Recognition*. New York: Prentice Hall.

图 模 型

图模型可视地表示变量之间的相互影响，并具有一个优点：利用条件独立性，可以将大量变量上的推断分解成一组涉及少量变量的局部计算。在给出一些手工推断的例子之后，我们讨论 d -分离和各种图上的信念传播算法。

16.1 引言

图模型 (graphical model) 又称贝叶斯网络 (Bayesian network)、信念网络 (belief network) 或概率网络 (probabilistic network)，它由节点和节点之间的有向弧构成。每个节点对应于一个随机变量 X ，并且具有一个对应于该随机变量的概率值 $P(X)$ 。如果存在一条从节点 X 到节点 Y 的有向弧，则表明 X 对 Y 有直接影响 (direct influence)，这一影响被条件概率 $P(Y|X)$ 所指定。网络是一个有向无环图 (directed acyclic graph, DAG)，即图中没有环。节点和节点之间的弧定义了网络的结构，而条件概率是给定结构的参数。

一个简单的例子在图 16-1 中给出，它对下雨 (R) 引起草地变湿 (W) 建模。天下雨的可能性为 40%，并且下雨时草地变湿的可能性为 90%；也许 10% 的时间雨下得不够，不足以让我们认为草地被淋湿。在这个例子中，随机变量是二元的，它们或者为真或者为假。存在 20% 的可能性草地变湿而实际上并没有下雨，例如，使用喷水器时。

我们可以看到三个值就可以完全指定联合分布 $P(R, W)$ 。如果 $P(R) = 0.4$ ，则 $P(\sim R) = 0.6$ 。类似地， $P(\sim W|R) = 0.1$ ，图 16-1 而对下雨是湿草地原因的贝叶斯网络建模而 $P(\sim W|\sim R) = 0.8$ 。联合概率表示成

$$P(R, W) = P(R)P(W|R)$$

通过在其父节点取所有可能值上求和，我们可以计算湿草地的 (边缘) 概率：

$$\begin{aligned} P(W) &= \sum_R P(R, W) = P(W|R)P(R) + P(W|\sim R)P(\sim R) \\ &= 0.9 \times 0.4 + 0.2 \times 0.6 = 0.48 \end{aligned}$$

如果我们知道下过雨，则湿草地的概率为 0.9；如果我们相信没有下过雨，则湿草地的概率低至 0.2；不知道是否下过雨，这个概率是 0.48。

图 16-1 显示了一个因果图 (causal graph)，解释草地变湿的主要原因是下雨。贝叶斯规则允许我们颠倒因果关系并且做出诊断 (diagnosis)。例如，已知草地是湿的，则下过雨的概率可以计算如下：

$$P(R|W) = \frac{P(W|R)P(R)}{P(W)} = 0.75$$

知道草地是湿的把下雨的概率由 0.4 增加到 0.75，这是因为 $P(W|R)$ 高，而 $P(W|\sim R)$ 低。我们通过添加节点和弧形成图，并产生依赖性。 X 和 Y 是独立事件 (independent event)，

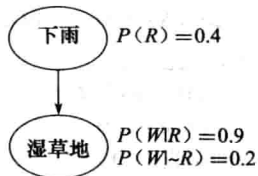


图 16-1 对下雨是湿草地原因的贝叶斯网络建模

如果

$$P(X, Y) = P(X)P(Y) \quad (16.1)$$

给定第三个事件 Z , X 和 Y 是条件独立事件 (conditional independent event), 如果

$$P(X, Y | Z) = P(X | Z)P(Y | Z) \quad (16.2)$$

这也可以写成

$$P(X | Y, Z) = P(X | Z) \quad (16.3)$$

在图模型中, 并非所有的节点都是连接的; 实际上, 一个节点一般只连接少数其他节点。特定的子图蕴含条件独立性陈述, 并且这些使得我们可以将一个复杂的图分解成较小的子集, 其中可以局部地做出推断, 并且稍后其结果在图上传播。有三种典型情况, 并且可以使用它们作为子图来构造较大的图。

16.2 条件独立的典型情况

1. 情况 1: 头到尾连接 (head-to-tail connection)

三个事件可以顺序连接, 如图 16-2(a) 所示。这里, 我们看到, 给定 Y , X 与 Z 是独立的: 知道 Y 就知道 Z 的一切, 知道 X 的状态并不能为 Z 增加附加信息, 记作 $P(Z | Y, X) = P(Z | Y)$ 。我们说 Y 阻塞 (block) 从 X 到 Z 的路径, 或者换句话说, Y 分离 (separate) X 和 Z , 意指如果删掉 Y , 则不存在 X 和 Z 之间的路径。在这种情况下, 联合概率写作

$$P(X, Y, Z) = P(X)P(Y | X)P(Z | Y) \quad (16.4)$$

这种表示联合概率的方式意味着独立性:

$$P(Z | X, Y) = \frac{P(X, Y, Z)}{P(X, Y)} = \frac{P(X)P(Y | X)P(Z | Y)}{P(X)P(Y | X)} = P(Z | Y) \quad (16.5)$$

典型地, X 是 Y 的原因, 而 Y 是 Z 的原因。例如, 如图 16-2(b) 所示, X 可以是多云 (C), Y 可以是下雨 (R), 而 Z 可以是湿草地 (W)。我们可以沿着链传播信息。如果我们不知道多云状态, 则有

$$P(R) = P(R | C)P(C) + P(R | \sim C)P(\sim C) = 0.38$$

$$P(W) = P(W | R)P(R) + P(W | \sim R)P(\sim R) = 0.48$$

假设早上我们看到天气是多云, 关于草地湿的概率我们能够说什么? 为此, 我们需要先将证据传播到中间节点 R , 然后传到询问节点 W 。

$$P(W | C) = P(W | R)P(R | C) + P(W | \sim R)P(\sim R | C) = 0.76$$

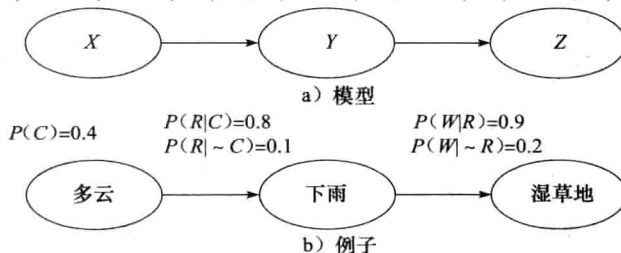


图 16-2 头到尾连接。a) 三个节点顺序连接。给定中间节点 Y , X 和 Z 是独立的: $P(Z | Y, X) = P(Z | Y)$ 。b) 例子, 天气多云导致下雨, 而下雨又导致湿草地

知道天气多云提高了湿草地的概率。我们还可以使用贝叶斯规则向后传播证据。假设我们旅游归来,看到我们的草地是湿的,那天多云的概率是多少?我们使用贝叶斯规则逆转方向:

$$P(C|W) = \frac{P(W|C)P(C)}{P(W)} = 0.65$$

知道草地湿将天气多云的概率从它的默认(先验)值 0.4 提高到 0.65。

2. 情况 2: 尾到尾连接(tail-to-tail connection)

X 可能是两个节点 Y 和 Z 的父节点,如图 16-3(a)所示。联合密度可以写作

$$P(X, Y, Z) = P(X)P(Y|X)P(Z|X) \quad (16.6)$$

通常, Y 和 Z 是通过 X 依赖的;给定 X , 它们变成独立的:

$$P(Y, Z|X) = \frac{P(X, Y, Z)}{P(X)} = \frac{P(X)P(Y|X)P(Z|X)}{P(X)} = P(Y|X)P(Z|X) \quad (16.7)$$

当它的值已知时, X 阻塞了 Y 和 Z 之间的路径,或换言之, X 分离了它们。

在图 16-3(b)中,我们看到一个例子,那里多云(C)天气影响下雨(R)和使用喷水器(S),一个是正的而另一个是负的。例如,知道下雨,我们可以使用贝叶斯规则逆转依赖性,并推断原因:

$$\begin{aligned} P(C|R) &= \frac{P(R|C)P(C)}{P(R)} = \frac{P(R|C)P(C)}{\sum_C P(R, C)} \\ &= \frac{P(R|C)P(C)}{P(R|C)P(C) + P(R|\sim C)P(\sim C)} = 0.89 \end{aligned} \quad (16.8)$$

注意,这个值大于 $P(C)$;知道下雨提高了天气多云的概率。

在图 16-3(a)中,例如,如果 X 未知但知道 Y ,则我们可以推断 X ,然后使用它推断 Z 。在图 16-3(b)中,知道喷水器的状态对下雨的概率有影响。如果我们知道喷水器在工作,则

$$\begin{aligned} P(R|S) &= \sum_C P(R, C|S) = P(R|C)P(C|S) + P(R|\sim C)P(\sim C|S) \\ &= P(R|C) \frac{P(S|C)P(C)}{P(S)} + P(R|\sim C) \frac{P(S|\sim C)P(\sim C)}{P(S)} = 0.22 \end{aligned} \quad (16.9)$$

这小于 $P(R) = 0.45$,即知道喷水器在工作降低了下雨的概率,因为喷水器和下雨发生在多云天气的不同状态。如果知道喷水器未工作,则用同样的方法我们发现 $P(R|\sim S) = 0.55$;这次,下雨的概率提高了。

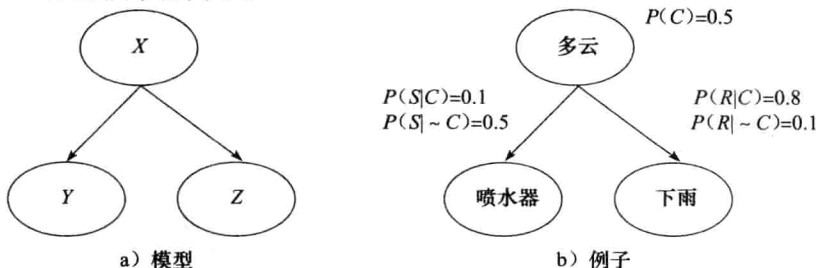


图 16-3 尾到尾连接。 X 是两个节点 Y 和 Z 的父节点。给定父节点,两个子节点是独立的: $P(Y|X, Z) = P(Y|X)$ 。在这个例子中,多云天气导致下雨,也使得我们不太可能使用喷水器

3. 情况3: 头到头连接(head-to-head connection)

在头到头节点中, 有两个父节点 X 和 Y 连接到单个节点 Z , 如图 16-4(a) 所示。联合密度记作

$$P(X, Y, Z) = P(X)P(Y)P(Z | X, Y) \quad (16.10)$$

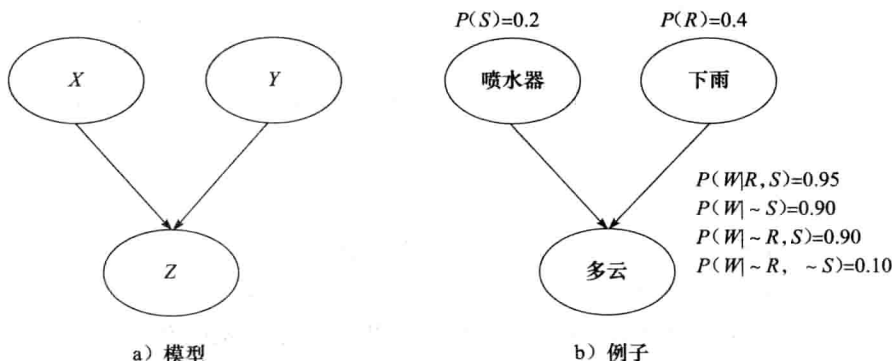


图 16-4 头到头连接。一个节点有两个父节点; 除非给定子节点, 否则两个父节点是独立的。

例如, 一个事件可能有两个独立的原因

X 和 Y 是独立的: $P(X, Y) = P(X)P(Y)$ (习题 2); 当已知 Z 时, 它们变成依赖的。对于这种情况, 阻塞或分离的概念不同: 当观察不到 Z 时, X 和 Y 之间的路径被阻塞, 或它们是分离的; 当 Z (或者它的任意后代) 被观测到时, 它们不再是被阻塞的、分离的, 也不是独立的。

392

例如, 在图 16-4(b) 中我们看到, 节点 W (湿草地) 有两个父节点: R (下雨) 和 S (喷水器), 因此它的概率是这两个值上的条件概率 $P(W | R, S)$ 。

不知道其他事情, 草地湿的概率通过联合概率边缘化来计算:

$$\begin{aligned} P(W) &= \sum_{R, S} P(W, R, S) \\ &= P(W | R, S)P(R, S) + P(W | \sim R, S)P(\sim R, S) \\ &\quad + P(W | R, \sim S)P(R, \sim S) + P(W | \sim R, \sim S)P(\sim R, \sim S) \\ &= P(W | R, S)P(R)P(S) + P(W | \sim R, S)P(\sim R)P(S) \\ &\quad + P(W | R, \sim S)P(R)P(\sim S) + P(W | \sim R, \sim S)P(\sim R)P(\sim S) = 0.52 \end{aligned}$$

现在, 假设我们知道喷水器开着, 则可以检查这对这个概率的影响。这是一个因果(预测)推理:

$$\begin{aligned} P(W | S) &= \sum_R P(W, R | S) \\ &= P(W | R, S)P(R | S) + P(W | \sim R, S)P(\sim R | S) \\ &= P(W | R, S)P(R) + P(W | \sim R, S)P(\sim R) = 0.92 \end{aligned}$$

我们看到 $P(W | S) > P(W)$, 知道喷水器工作提高了湿草地的概率。

给定草地是湿的, 我们也可以计算喷水器开着的概率。这是诊断推理。

$$P(S | W) = \frac{P(W | S)P(S)}{P(W)} = 0.35$$

$P(S|W) > P(S)$, 即知道草地是湿的提高了喷水器开着的概率。现在让我们假设下过雨, 于是, 我们有

$$P(S|R,W) = \frac{P(W|R,S)P(S|R)}{P(W|R)} = \frac{P(W|R,S)P(S)}{P(W|R)} = 0.21$$

这个值比 $P(S|W)$ 小, 称作解释远离(explaining away); 给定我们已知下过雨, 则喷水器导致湿草地的概率降低了。已知草地是湿的, 下雨和喷水器成为相互依赖的。类似地, $P(S|\sim R, W) > P(S|W)$ 。当我们比较 $P(R|W)$ 和 $P(R|W, S)$ 的时候, 看到类似的情况(习题 3)。

我们可以通过组合这样的子图来构造更大的图。例如, 在图 16-5 中, 我们组合了前面的两个子图, 可以计算如果多云草地湿的概率:

$$\begin{aligned} P(W|C) &= \sum_{R,S} P(W,R,S|C) \\ &= P(W,R,S|C) + P(W,\sim R,S|C) + P(W,R,\sim S|C) + P(W,\sim R,\sim S|C) \\ &= P(W|R,S,C)P(R,S|C) + P(W|\sim R,S,C)P(\sim R,S|C) \\ &\quad + P(W|R,\sim S,C)P(R,\sim S|C) + P(W|\sim R,\sim S,C)P(\sim R,\sim S|C) \\ &= P(W|R,S)P(R|C)P(S|C) + P(W|\sim R,S)P(\sim R|C)P(S|C) \\ &\quad + P(W|R,\sim S)P(R|C)P(\sim S|C) + P(W|\sim R,\sim S)P(\sim R|C)P(\sim S|C) \end{aligned}$$

其中, 我们使用了 $P(W|R,S,C) = P(W|R,S)$; 给定 R 和 S , W 独立于 C : R 和 S 阻塞了 W 和 C 之间的路径。类似地, $P(R,S|C) = P(R|C)P(S|C)$; 给定 C , R 和 S 是独立的。这里, 我们看到贝叶斯网络的优点: 它明确地表示了独立性, 并且使得我们能够将推断分解成若干从证据节点到查询节点传播小的变量组上的计算。

我们可以计算 $P(C|W)$, 并且有诊断推理:

$$P(C|W) = \frac{P(W|C)P(C)}{P(W)}$$

图形表示是可视的并且有助于理解。这个网络提供了条件独立性陈述, 并且允许我们将多个变量的联合分布问题分解成局部结构; 这简化了分析和计算。图 16-5 表示了一个四个二元变量的联合密度, 它通常需要存储 15 个值($2^4 - 1$), 而这里只有 9 个。如果每个节点只有少量的父节点, 则复杂度由指数降到线性(按节点数)。正如我们在前面看到的, 当联合密度分解成较小变量组的条件密度时, 推断也会变得更容易:

$$P(C,S,R,W) = P(C)P(S|C)P(R|C)P(W|S,R) \quad (16.11)$$

在一般情况下, 当我们有变量 X_1, \dots, X_d 时, 有

$$P(X_1, \dots, X_d) = \prod_{i=1}^d P(X_i | \text{parents}(X_i)) \quad (16.12)$$

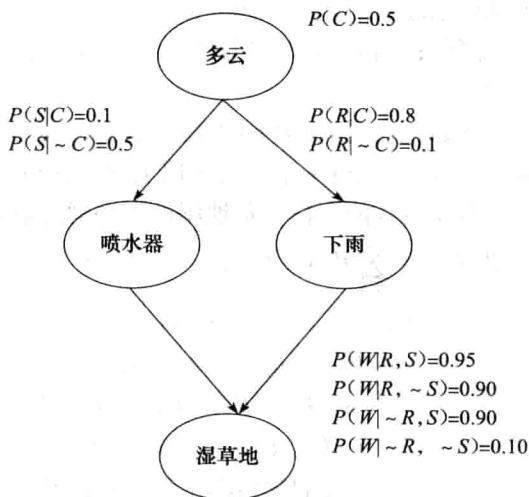


图 16-5 通过组合较简单的子图形成较大的图, 使用隐含的条件独立性, 使信息在图上传播

于是, 给定 X_i 的任意子集, 即根据证据赋予它们一定的值, 通过边缘化联合分布, 我们可以计算 X_i 的其他子集的概率分布。这开销很大, 因为它需要计算指数个联合概率组合, 即使每个都能像式 (16.11) 那样被简化。然而, 给定相同的证据, 对于不同的 X_i , 我们可以使用相同的中间值(条件概率的乘积和边缘化的和), 并且在 16.5 节我们将讨论信念传播算法, 通过进行一次可以用于不同查询节点的多次局部中间计算, 以更小开销做推断。

尽管在这个例子中我们使用了二元变量, 但是可以将它直接推广到具有任意多个可能值的离散变量(具有 m 个可能值和 k 个父节点, 需要大小为 m^k 的条件概率表)或连续变量(参数化的, 例如 $p(Y|x) \sim \mathcal{N}(\mu(x|\theta), \sigma^2)$, 参见 16.3.3 节)。

使用贝叶斯网络的一个主要优点是我们不必明确指定某些变量作为输入, 某些其他变量作为输出。任何变量集的值都能通过证据建立, 而任何其他变量集的概率都可以推断, 并且非监督学习和监督学习之间的区别变得模糊不清。从这个角度看, 一个图模型可以看作一个“概率数据库”(Jordan 2009), 一种可以回答关于随机变量值的查询机器。

在一个问题中, 还可能隐藏变量(hidden variable), 它们的值不能通过证据获得。使用隐藏变量的优点是可以更容易地定义依赖结构。例如, 在购物篮分析中, 当我们想找出销售的商品之间的依赖性时, 比如我们知道“婴儿食品”、“尿布”和“牛奶”之间的依赖性, 即顾客购买其中一种, 多半可能购买另外两种, 我们不是将(非因果的)弧放在这三者之间, 而可能是指定一个隐藏节点“家庭婴儿”作为这三种商品消费的隐藏原因。当存在隐藏节点时, 它们的值用观测节点的给定值估计并填写。

需要强调一点, 从节点 X 到 Y 的链不是也不必总是意味着因果关系(causality)。它只是意味着 X 在 Y 上有直接影响(direct influence), 即 Y 的概率以 X 的值为条件, 并且即使没有直接的原因, 两个节点之间仍可能有一个边。通过提供关于数据如何产生的解释, 在构造网络时有因果关系更可取(Pearl 2000), 但是这种因果关系并非总是能够得到的。

16.3 图模型实例

16.3.1 朴素贝叶斯分类

对于分类, 对应的图模型如图 16-6(a) 所示, 其中 \mathbf{x} 是输入, C 是一个多元变量, 取类编码的 K 个状态之一。正如我们在图 16-1 的下雨和湿草地的例子中看到的, 贝叶斯规则允许诊断:

$$P(C|\mathbf{x}) = \frac{P(C)p(\mathbf{x}|C)}{P(\mathbf{x})}$$

如果输入是相互独立的, 则我们有如图 16-6(b) 所示的图, 它称作朴素贝叶斯分类器(naive Bayes' classifier), 因为它忽略了输入之间可能的依赖性(即相关性), 将一个多变量问题归约为一组单变量问题:

$$p(\mathbf{x}|C) = \prod_{j=1}^d p(x_j|C)$$

在 5.5 和 5.7 节, 我们已经分别对数值和离散的 \mathbf{x} 讨论过这种情况的分类。

聚类也类似, 不同之处在于多项式类指示变量 C 在分类中是可观测的, 而类似的簇指

示变量 Z 是不可观测的。期望最大化算法(7.4 节)的 E 步使用贝叶斯规则逆转弧的方向, 并估计给定输入的簇指示子。

图 16-6(a) 是一个创建数据过程的生成模型(generative model)。仿佛我们先通过由 $P(C)$ 抽样, 随机地选择类 C , 然后固定 C , 通过从 $p(\mathbf{x} | C)$ 抽样选取 \mathbf{x} 。将数据看作由一个可以可视化为图的因果生成模型抽样不仅易于理解, 而且也便于在许多领域进行推断。例如, 在文本分类中产生一个文本可以看作这样一个过程: 一位作者决定写一篇关于某特定主题的文章, 然后选择相应的词集合。在生物信息学中, 众多使用图形方法的领域之一是系统演化树(phylogenetic tree)。系统演化树是一个有向图, 其树叶是当前物种, 非终端节点是过去的祖先, 它们在物种形成事件中进化成多个物种, 而条件概率依赖于物种和它的祖先之间的进化距离(Jordan 2004)。

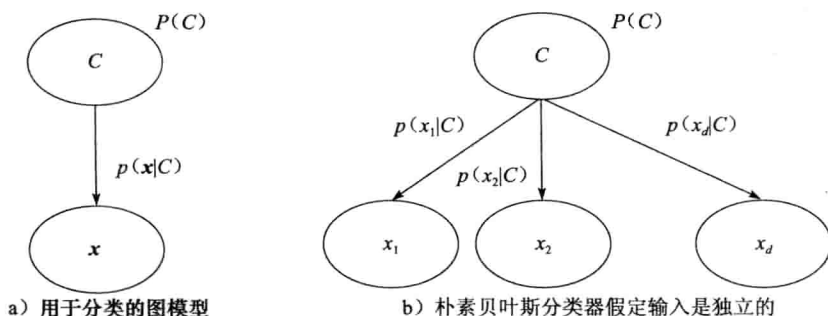


图 16-6

16.3.2 隐马尔可夫模型

隐马尔可夫模型(hidden Markov model, HMM)在前面的第 15 章讨论过, 它是情况 1 的一个例子, 其中相继的节点 q_{t-2} , q_{t-1} , q_t 对应于一阶马尔可夫模型中链上的三个状态。时刻 t 的状态 q_t 仅依赖于时刻 $t-1$ 时的状态 q_{t-1} , 并且给定 q_{t-1} , q_t 独立于 q_{t-2}

$$P(q_t | q_{t-1}, q_{t-2}) = P(q_t | q_{t-1})$$

正如状态转移概率矩阵 A 所给定的(参见图 16-7)。每个隐藏变量产生一个被观测的离散观测, 正如观测概率矩阵 B 给定的。隐马尔可夫模型的前后向过程是我们马上将要讨论的信念传播的特例。

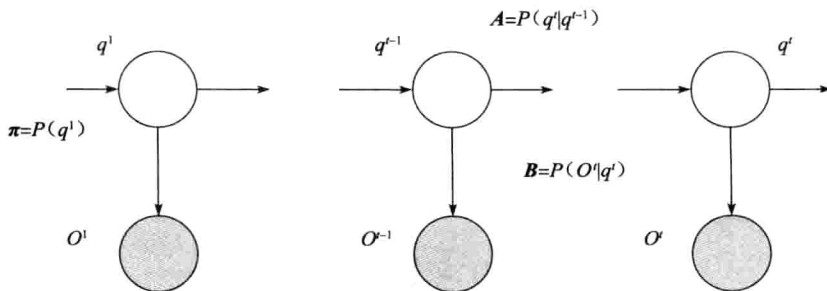


图 16-7 隐马尔可夫模型可以绘制成图模型, 其中 q^i 是隐藏状态, 带阴影的 O^i 是被观测的

不同类型的 HMM 可以用不同的图模型表示。在图 16-8(a)中,显示了一个输入-输出 HMM(input-output HMM)(参见 15.9 节),其中有两个分离的观测输入-输出序列,还有一个隐藏状态序列。输出观测依赖于状态,也依赖于输入;可以将这看作矩阵 B , 其元素不是标量,而是输入的参数化函数。这也可以类似地看作混合专家架构(参见 12.8 节),其选通输出(隐藏状态)还依赖于前一时间步的选通值。

另一类容易可视化的 HMM 是因子 HMM(factorial HMM),其中有多多个分离的隐藏序列,它们相互影响,产生单个观测序列。一个例子是显示父母-子女关系的家谱(pedigree)(Jordan 2004);图 16-8(b)对减数分裂(meiosis)建模,其中两个序列对应于父亲和母亲(他们是独立的)的染色体,并且在每个位点(基因),后代从父亲那里接受一个等位基因,或从母亲那里接受其他等位基因。

图 16-8(c)显示的耦合 HMM(coupled HMM)对产生两个平行观测序列的两个平行但相关的隐藏序列建模。例如,在语音识别中,我们可以有一个观测的讲话声学序列和一个观测的唇动图像的可视序列,每个序列都有自己的隐藏状态,隐藏状态中的二者是依赖的。

在图 16-8(d)中显示的开关 HMM(switching HMM)中,有 K 个平行的、独立的隐藏状态序列,而状态变量 S 在任何时刻都挑选其中一个序列,并且选中的序列产生输出。也就是说,我们随着进展在状态序列之间切换。

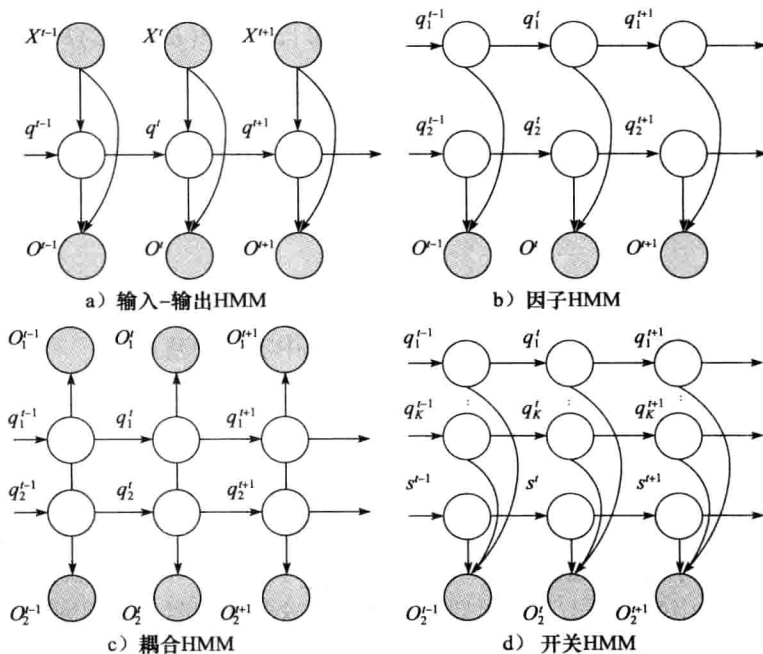


图 16-8 不同的 HMM 模型,关于从潜在变量的马尔可夫序列产生的观测数据(用阴影显示)方法的不同假设

在典型的 HMM 中,尽管观测可能是连续的,但是状态是离散的。线性动态系统(linear dynamical system)又称卡尔曼滤波器(Kalman filter),其中状态和观测都是连续的。在基本情况下,时刻 t 状态是时刻 $t-1$ 状态的具有零均值高斯加性噪声的线性函数,并且在每个状态中,观测是另一个具有零均值高斯加性噪声的线性函数。这两个线性映射和两个噪声源的协方差构成参数。我们先前讨论的任何 HMM 的变体都可以类似地推广到使用连续状态。

通过适当地修改图模型,可以使该结构适应产生数据过程的特点。这种将模型与数据匹配的过程是权衡偏倚和方差的最佳模型选择过程。缺点是在这种扩充的 HMM 上,精确推断已经不再可能,并且需要近似或抽样方法 (Ghahramani 2001, Jordan 2009)。

16.3.3 线性回归

线性回归可以看作一个图模型,如图 16-9 所示。输入 \mathbf{x}' 从先验 $p(\mathbf{x})$ 中抽取,而因变量 r' 依赖于输入 \mathbf{x} 、权重 \mathbf{w} (从被 α 参数化的先验抽取,即从 $p(\mathbf{w}) \sim \mathcal{N}(\mathbf{0}, \alpha^{-1}\mathbf{I})$ 中抽取) 和噪声 ε (被 β 参数化,即 $p(\varepsilon) \sim \mathcal{N}(0, \beta^{-1})$):

$$p(r' | \mathbf{x}', \mathbf{w}) \sim \mathcal{N}(\mathbf{w}^T \mathbf{x}', \beta^{-1}) \quad (16.13)$$

训练集中有 N 个这样的对,显示在图 16-9 的矩形框中。给定一个新的输入 \mathbf{x}' ,目标是估计 r' ,它将是 $E[r' | \mathbf{x}', \mathbf{w}]$ 。

权重 \mathbf{w} 未给出,但可以使用训练集 $[\mathbf{X}, \mathbf{r}]$ 估计它们。在式(16.9)中, C 是 R 和 S 的原因,其中我们用

$$P(R | S) = \sum_C P(R, C | S) = P(R | C)P(C | S) + P(R | \sim C)P(\sim C | S)$$

用 S 替代 C ,而 S 又用来估计 R 。类似地,这里我们有

$$\begin{aligned} p(r' | \mathbf{x}', \mathbf{r}, \mathbf{X}) &= \int p(r' | \mathbf{x}', \mathbf{w}) p(\mathbf{w} | \mathbf{X}, \mathbf{r}) d\mathbf{w} \\ &= \int p(r' | \mathbf{x}', \mathbf{w}) \frac{p(\mathbf{r} | \mathbf{X}, \mathbf{w}) p(\mathbf{w})}{p(\mathbf{r})} d\mathbf{w} \\ &\propto \int p(r' | \mathbf{x}', \mathbf{w}) \prod_i p(r_i | \mathbf{x}_i, \mathbf{w}) p(\mathbf{w}) d\mathbf{w} \end{aligned} \quad (16.14)$$

其中,第二行根据贝叶斯规则,而第三行根据训练集中实例的独立性。

16.4 d -分离

现在,我们用 d -分离(d -separation)推广阻塞和分离的概念,并用这样的方式定义它,使得对于节点的任意子集 A 、 B 和 C ,我们可以检查给定 C , A 和 B 是否是独立的。Jordan (2009)将这想象成一个球在图上跳动,并称为贝叶斯球(Bayes's ball)。我们将 C 中的节点设置为它们的值,在 A 中的每个节点上放置一个球,让这些球按照一组规则四处移动,并检查是否有一个球到达 B 中某个节点。如果是,则它们是依赖的;否则,它们是独立的。

为了检查给定 C , A 和 B 是否是 d -分离的,我们考虑 A 中任意节点与 B 中任意节点之间的所有可能路径。任意一条这样的路径是阻塞的(blocked),如果

- 1) 路径上边的方向或者满足头到尾(情况 1),或者满足尾到尾(情况 2),并且该节点在 C 中;或者
- 2) 路径上边的方向满足头到头(情况 3),并且无论哪个节点还是它的任意后代都不在 C 中。

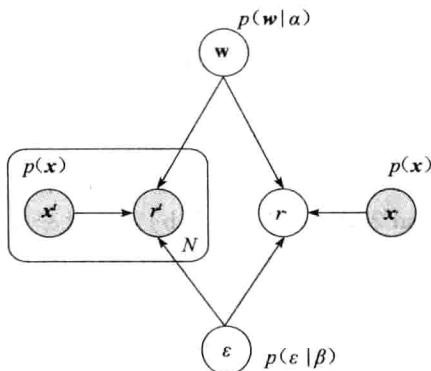


图 16-9 线性回归的贝叶斯网络

如果所有的路径都是阻塞的，则认为 A 和 B 是 d -分离的，即给定 C 它们是独立的；否则，它们是依赖的，例子在图 16-10 中。

16.5 信念传播

我们已经讨论了一些手工推断的例子。现在，我们感兴趣的是回答诸如 $P(X|E)$ 这种查询的算法，其中 X 是图中的任意查询节点(query node)，而 E 是其值已被设置为确定值的证据节点(evidence node)的任意子集。按照 Pearl(1988) 的做法，我们从链这种最简单的情况开始，逐渐地考虑更复杂的图。我们的目标是找到诸如贝叶斯规则或边缘化这样的概率过程的对应图操作，使得推断任务可以映射到通用的图算法。

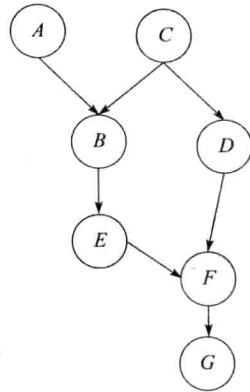


图 16-10 d -分离的例子。给定 C ，路径 $BCDF$ 是被阻塞的，因为 C 是一个尾到尾节点。 $BEFG$ 被 F 阻塞，因为 F 是一个头到尾节点。 $BEFD$ 是被阻塞的，除非给定 F (或 G)

16.5.1 链

链(chain)是头到尾节点的序列，有一个没有父节点的根(root)节点，其他所有节点都只有一个父节点；除最后一个称作叶(leaf)的节点外，其他所有节点都有一个子节点。如果证据在 X 的祖先中，则我们只能做诊断推断，并沿着链向上传播证据；如果证据在 X 的后代中，则我们可以使用贝叶斯规则做因果推断并向上传播。让我们看一般情况：我们在两个方向上都有证据，上链 E^+ 和下链 E^- (参见图 16-11)。注意，任何证据节点都将 X 与链上证据另一侧的节点分离，并且它们的值不影响 $P(X)$ ；这对两个方向都成立。

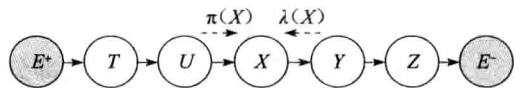


图 16-11 沿链推断

我们把每个节点都看作一个处理机，它从它的邻居接收消息并在局部计算后沿着链传递。每个节点 X 都局部地计算并存储两个值： $\lambda(X) \equiv P(E^- | X)$ 是传播的 E^- ，其中 X 从它的子节点接收并传给它的父节点；而 $\pi(X) \equiv P(E^+ | X)$ 是传播的 E^+ ，其中 X 从它的父节点接收并传给它的子节点。

$$\begin{aligned}
 P(X|E) &= \frac{P(E|X)P(X)}{P(E)} = \frac{P(E^+, E^- | X)P(X)}{P(E)} \\
 &= \frac{P(E^+ | X)P(E^- | X)P(X)}{P(E)} \\
 &= \frac{P(X|E^+)P(E^+)P(E^- | X)P(X)}{P(X)P(E)} \\
 &= \alpha P(X|E^+)P(E^- | X) = \alpha \pi(X)\lambda(X) \quad (16.15)
 \end{aligned}$$

对于某个不依赖于 X 值的规范化常量 α ，第二行是因为给定 X ， E^+ 和 E^- 是独立的，而第三行是根据贝叶斯规则。

如果节点 E 被例示为确定的值 \tilde{e} ，则 $\lambda(\tilde{e}) \equiv 1$ ，而对于 $e \neq \tilde{e}$ 有 $\lambda(e) \equiv 0$ 。没有例示的

叶节点 X 对所有的 x 值有 $\lambda(x) \equiv 1$ 。没有例示的根节点 X 取先验概率作为 π 值：对于任意 x , $\pi(x) \equiv P(x)$ 。

给定这些初始条件, 我们可以设计一个沿着链传播证据的递归方案。

对于 π -消息, 我们有

$$\begin{aligned}\pi(X) &\equiv P(X|E^+) = \sum_U P(X|U, E^+) P(U|E^+) \\ &= \sum_U P(X|U) P(U|E^+) = \sum_U P(X|U) \pi(U)\end{aligned}\quad (16.16)$$

其中, 第二行基于 U 阻塞 X 与 E^+ 之间的路径这一事实。

对于 λ -消息, 我们有

$$\begin{aligned}\lambda(X) &\equiv P(E^-|X) = \sum_Y P(E^-|X, Y) P(Y|X) \\ &= \sum_Y P(E^-|Y) P(Y|X) = \sum_Y P(Y|X) \lambda(Y)\end{aligned}\quad (16.17)$$

其中, 第二行基于 Y 阻塞 X 与 E^- 之间的路径这一事实。

当证据节点的值被设置时, 它们启动传播, 并且节点持续更新, 直至收敛。Pearl (1988) 将这看作一台并行机器, 其中每个节点用一个处理器实现, 每个处理器都与其他处理器并行, 通过 π -消息和 λ -消息与它的父节点和子节点交换信息。

16.5.2 树

链是受限的, 因为每个节点只能有一个父节点和一个子节点, 即单个原因和单个症状。在树中, 每个节点都可以有多个子女, 但是除单个根外, 所有的节点都只有一个父节点。相同的信念传播也可以在树上进行, 与链的不同之处是: 节点由它的诸子女接收不同的 λ -消息, 并且向它的诸子女发送不同的 π -消息。 $\lambda_Y(X)$ 表示 X 从它的子女 Y 接收的消息, $\pi_Y(X)$ 表示 X 发送到它的子女 Y 的消息。

我们还是将可能的证据划分成两部分: E^- 是以查询节点 X 为根的子树中的证据节点, 而 E^+ 是其他证据节点 (参见图 16-12)。注意, E^+ 中的节点不必是 X 的祖先, 而可以在以 X 的兄妹节点为根的子树中。重要的是, X 分离了 E^+ 和 E^- , 使得我们可以有 $P(E^+, E^-|X) = P(E^+|X)P(E^-|X)$, 因而有

$$P(X|E) = \alpha \pi(X) \lambda(X)$$

其中, α 是一个规范化常量。

$\lambda(X)$ 是以 X 为根的子树中的证据, 并且如图 16-12 所示, 如果 X 有两个子女 Y 和 Z , 则它用下式计算:

$$\begin{aligned}\lambda(X) &\equiv P(E_X^-|X) = P(E_Y^-, E_Z^-|X) \\ &= P(E_Y^-|X) P(E_Z^-|X) = \lambda_Y(X) \lambda_Z(X)\end{aligned}\quad (16.18)$$

在一般情况下, 如果 X 有 m 个子女 $Y_j (j=1, \dots, m)$, 则我们将它们的所有 λ 值相乘:

$$\lambda(X) = \prod_{j=1}^m \lambda_{Y_j}(X) \quad (16.19)$$

一旦 X 从它的子女的 λ -消息中收集了 λ 证据, 它就将这些证据上传给它的父节点:

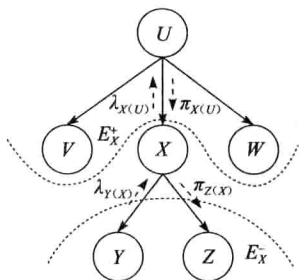


图 16-12 在树中, 一个节点可能有多个子女, 但只有一个父节点

$$\lambda_X(U) = \sum_X \lambda(X) P(X|U) \quad (16.20)$$

类似地, 在另一个方向, $\pi(X)$ 是收集在 $P(U|E^+)$ 中其他地方的证据, 并且作为 π -消息传到 X :

$$\pi(X) \equiv P(X|E_X^+) = \sum_U P(X|U) P(U|E_X^+) = \sum_U P(X|U) \pi_X(U) \quad (16.21)$$

然后, 这个计算的 π 值下传到 X 的子女。注意, Y 从 X 接收的是 X 从它的父节点 U 和它的其他子女 Z 接收的, 它们一起组成 E_Y^+ (参见图 16-12):

$$\begin{aligned} \pi_Y(X) \equiv P(X|E_Y^+) &= P(X|E_X^+, E_Z^-) = \frac{P(E_Z^-|X, E_X^+) P(X|E_X^+)}{P(E_Z^-)} = \frac{P(E_Z^-|X) P(X|E_X^+)}{P(E_Z^-)} \\ &= \alpha \lambda_Z(X) \pi(X) \end{aligned} \quad (16.22)$$

同样, 如果 Y 不止 Z 一个兄妹, 而是有多个, 则需要取它们所有的 λ 值的乘积:

$$\pi_{Y_j}(X) = \alpha \prod_{i \neq j} \lambda_{Y_i}(X) \pi(X) \quad (16.23)$$

16.5.3 多树

在树中, 节点有单个父节点, 即单个原因。在多树 (polytree) 中, 一个节点可以有多个父节点, 但是我们要求图是单连接的, 这意味着任意两个节点之间只有一条链。如果我们删除 X , 则图就被划分成两个分支。这是必要的, 使得我们可以继续将 E_X 划分成 E_X^+ 和 E_X^- , 给定 X , 则它们是独立的 (参见图 16-13)。

如果 X 有多个父节点 $U_i (i=1, \dots, k)$, 则它从所有父节点接收 π -消息 $\pi_X(U_i)$, 按如下方法组合它们:

$$\begin{aligned} \pi(X) \equiv P(X|E_X^+) &= P(X, E_{U_1X}^+, E_{U_2X}^+, \dots, E_{U_kX}^+) \\ &= \sum_{U_1} \sum_{U_2} \dots \sum_{U_k} P(X|U_1, U_2, \dots, U_k) \\ &\quad P(U_1|E_{U_1X}^+) \dots P(U_k|E_{U_kX}^+) \\ &= \sum_{U_1} \sum_{U_2} \dots \sum_{U_k} P(X|U_1, U_2, \dots, U_k) \prod_{i=1}^k \pi_X(U_i) \end{aligned} \quad (16.24)$$

并将它传递到多个子女 $Y_j (j=1, \dots, m)$:

$$\pi_{Y_j}(X) = \alpha \prod_{i \neq j} \lambda_{Y_i}(X) \pi(X) \quad (16.25)$$

在 X 具有多个父节点时, X 传递到它其中的一个父节点 U_i 的 λ -消息不仅组合了 X 从它的子女接收的证据, 而且还有 X 从它的其他父节点 $U_r (r \neq i)$ 接收的 π -消息, 它们一起构成 $E_{U_iX}^-$:

$$\begin{aligned} \lambda_X(U_i) \equiv P(E_{U_iX}^-|X) &= \sum_X \sum_{U_{r \neq i}} P(E_X^-, E_{U_{r \neq i}X}^+, X, U_{r \neq i} | U_i) \\ &= \sum_X \sum_{U_{r \neq i}} P(E_X^-, E_{U_{r \neq i}X}^+ | X, U_{r \neq i}, U_i) P(X, U_{r \neq i} | U_i) \\ &= \sum_X \sum_{U_{r \neq i}} P(E_X^- | X) P(E_{U_{r \neq i}X}^+ | U_{r \neq i}) P(X | U_{r \neq i}, U_i) P(U_{r \neq i} | U_i) \end{aligned}$$

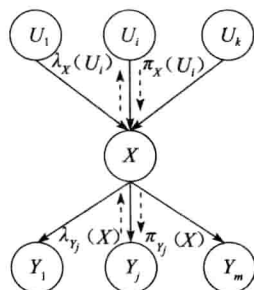


图 16-13 在多树中, 一个节点可以有多个子节点和多个父节点, 但是图是单连接的, 即 U_i 和 Y_j 之间只有一条经过 X 的链

$$\begin{aligned}
&= \sum_X \sum_{U_{r \neq i}} P(E_X^- | X) \frac{P(U_{r \neq i} | E_{U_{r \neq i} X}^+) P(E_{U_{r \neq i} X}^+)}{P(U_{r \neq i})} P(X | U_{r \neq i}, U_i) P(U_{r \neq i} | U_i) \\
&= \beta \sum_X \sum_{U_{r \neq i}} P(E_X^- | X) P(U_{r \neq i} | E_{U_{r \neq i} X}^+) P(X | U_{r \neq i}, U_i) \\
&= \beta \sum_X \sum_{U_{r \neq i}} \lambda(X) \prod_{r \neq i} \pi_X(U_r) P(X | U_1, \dots, U_k) \\
&= \beta \sum_X \lambda(X) \sum_{U_{r \neq i}} P(X | U_1, \dots, U_k) \prod_{r \neq i} \pi_X(U_r) \quad (16.26)
\end{aligned}$$

与树一样, 为了找出它的总 λ , 父节点取它从子女接收的 λ -消息的乘积:

$$\lambda(X) = \prod_{j=1}^m \lambda_{V_j}(X) \quad (16.27)$$

在这种有多个父节点的情况下, 我们需要存储和处理给定所有父节点的条件概率 $p(X | U_1, \dots, U_k)$, 对于大的 k , 这开销很大。目前已经提出了一些方法, 将复杂度从 k 的指数降低到线性。例如, 在噪声或门 (noisy OR gate) 中, 当多个父事件出现时, 任意父节点都足以导致该事件并且似然不减少。如果在仅有一个原因 (U_i) 出现时, X 发生的概率是 $1 - q_i$

$$P(X | U_i, \sim U_{p \neq j}) = 1 - q_i \quad (16.28)$$

当它们的一个子集 T 出现时, X 发生的概率用下式计算:

$$P(X | T) = 1 - \prod_{u_i \in T} q_i \quad (16.29)$$

例如, 假设湿草地有两个原因: 下雨和使用喷水器, 其概率为 $q_R = q_S = 0.1$, 即二者都有 90% 的可能性导致湿草地。于是, $P(W | R, \sim S) = 0.9$, $P(W | R, S) = 0.99$ 。

另一种可能是, 给定参数集, 将该条件概率写成某个函数, 例如一个线性模型

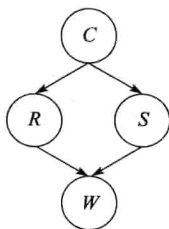
$$P(X | U_1, \dots, U_k, w_0, w_1, \dots, w_k) = \text{sigmoid} \left(\sum_{i=1}^k w_i U_i + w_0 \right) \quad (16.30)$$

其中, sigmoid 确保输出是 0 和 1 之间的概率。例如, 在训练阶段, 我们可以学习参数 $w_i (i = 1, \dots, d)$, 最大化样本上的最大似然。

16.5.4 结树

如果有环, 即如果基本无向图有环 (例如, 如果 X 的父母有共同的祖先), 则我们先前讨论的算法就不能运行。在这种情况下, 存在多条传播证据的路径, 例如, 在计算 X 上的概率时, 我们不能说 X 把 E 分离成 E_X^+ 和 E_X^- , 而分别作为原因 (向上) 和诊断 (向下) 证据; 删除 X 不会把图一分为二。在 X 上对它们取条件概率不会使它们独立, 而它们两个可能通过不经过 X 的其他路径相互影响。

如果可以将该图转换成多树, 则我们仍然可以使用同样的算法。我们定义团节点 (clique node) 对应于原变量的一个子集, 并且连接它们, 使得它们形成一棵树 (参见图 16-14)。于是, 利用这种改动, 我们可以运



a) 一个多连接图



b) 节点被聚类的对应结树

图 16-14

行同样的信念传播算法，这就是结树算法(junction tree algorithm)的基本思想(Lauritzen 和 Spiegelhalter 1988; Jensen 1996; Jordan 2009)。

16.6 无向图：马尔可夫随机场

迄今为止，我们讨论了有向图，其中影响是无向的，并且我们使用贝叶斯规则来逆转弧的方向。如果影响是对称的，则我们使用无向图模型来表示它们，这种模型也称作马尔可夫随机场(Markov random field)。例如，图像中的邻近像素趋向于具有相同的颜色(即相关的)，并且相关性沿着两个方向。

有向和无向图以不同的方式定义条件独立性，因而被有向图表示的概率分布不同于被无向图表示的概率分布，反之亦然(Pearl 1988)。

由于没有方向，因此没有弧的头尾之分，所以无向图的处理更简单。例如，给定 C ，检查 A 和 B 是否独立更加简单。我们只需检查如果删除 C 中的所有节点，我们是否还有一条从 A 中的一个节点到 B 中的一个节点的路径。如果有，则它们是依赖的。否则，如果 A 中节点与 B 中节点之间的所有路径都经由 C 中的节点，使得 C 的删除导致 A 中的节点和 B 中的节点在不同的分支中，则我们有独立性。

在无向图的情况下，我们不提及父节点或子女节点，而是讨论团(clique)，团是节点的集合，使得该集合中的任意两个节点之间都存在一条边。极大团(maximal clique)是元素数取极大的团。不采用条件概率(意味着方向)，在无向图中，我们有势函数(potential function) $\psi_c(X_c)$ ，其中 X_c 是团 C 中变量的集合，并且我们定义联合分布为图中极大团的势函数乘积：

$$p(X) = \frac{1}{Z} \prod_c \psi_c(X_c) \quad (16.31)$$

其中 Z 是正则化常数，以确保 $\sum_X p(X) = 1$ ：

$$Z = \sum_X \prod_c \psi_c(X) \quad (16.32)$$

可以证明，有向图已经是正则化的(习题5)。

与有向图不同，无向图中的势函数不必有概率解释，并且在定义它们时可以有更大的自由度。一般地，我们可以把势函数看作表现局部约束，即偏爱某种局部配置而不是其他局部配置。例如，在图像中，我们可以定义邻近像素之间的两两势函数，它们的颜色相似时的取值比颜色不同时的取值更大(Bishop 2006)。于是，设置某些像素的值作为证据，我们可以估计其他未知像素的值，例如，基于咬合(occlusion)。

如果我们有有向图，则容易通过丢弃所有方向将它重新绘制成无向图，并且如果一个节点只有一个父节点，则我们可以简单地令两两势函数转为条件概率。然而，如果一个节点有多个父节点，则由于头到头节点的“解释远离”现象使得这些父节点是依赖的，因而我们应该将这些父节点放在相同的团中，使得团的势包含所有的父节点。这通过如下方法来实现：用链连接节点的所有父节点，使得它们之间完全连接并形成团，这称作“嫁娶”(marrying)父母，而这一过程称作教化(moralization)。顺便说一下，教化是产生无向结树的步骤之一。

将信念传播算法用于无向图是直截了当的，并且更容易，因为势函数是对称的，并且不需要区分因果和诊断证据。这样，我们可以在无向链和无向树上做推断。但是在节点具有多

个父节点的多树中,教化必然导致环,因此这种方法不行。一个技巧是将其转换成因素图(factor graph),除了变量节点之外,使用第二种类型的因素节点(factor node),并且将联合分布表示成因素的乘积(Kschischang、Frey 和 Loeliger 2001):

$$p(X) = \frac{1}{Z} \prod_S f_S(X_S) \quad (16.33)$$

其中, X_S 代表被因素 S 使用的变量节点的子集。有向图是一个特例,其中因素对应于局部条件分布;无向图是另一个特例,其中因素是极大团上的势函数。正如我们在图 16-15 中所看到的,这样做的优点是教化后仍然能够保持树结构。

可以将信念传播算法推广到因素图,这称作和-积算法(sum-product algorithm)(Bishop 2006, Jordan 2009)。该算法具有同样的思想:做一次局部计算,并作为消息通过图来传播它们。不同的是,这里有两类消息,因为有因素和变量两类节点,所以我们要区分它们的消息。注意,因素图是一个二部图,一类节点只能与另一类节点直接连接。

在信念传播中或在和-积算法中,给定固定为某特定值的证据节点的集合 E ,目标是找出节

点集合 X 的概率,即求 $P(X|E)$ 。在某些应用中,我们感兴趣的可能是找出最大化联合概率分布 $p(X)$ 的所有 X 的设置。例如,在无向的情况下,势函数对局部一致性编码,这种方法将在整个图上传播局部一致性,并找出最大化全局一致性的解。在节点对应于像素而两两势函数支持相关性的图中,这种方法将实现噪声删除(Bishop 2006)。实现这种方法的算法称作最大-积算法(max-product algorithm)(Bishop 2006, Jordan 2009),它与和-积算法一样,但是取最大值(最可能的值),而不是取和(边缘化)。这类似于我们在第 15 章讨论的隐马尔可夫模型的正反向算法与 Viterbi 算法之间的区别。

注意,节点不必对应于诸如像素这样的低层概念。例如,在视频应用中,我们可以有针对性地对不同类型的角或不同方向线的节点,使用势函数检查相容性,以便观察它们是否可能是相同解释的一部分(例如,回忆 Necher 立方体),使得总体相容的解在组合局部证据之后显现。

推断算法在多树或结树上的复杂度由父节点的最大个数或最大团的大小决定,当这些很大时,严格推断是不可行的。在这种情况下,我们需要使用近似或抽样算法(Jordan 1999, Bishop 2006, Jordan 2009)。

16.7 学习图模型的结构

像任何学习方法一样,学习图模型有两部分工作。第一部分是给定结构学习参数。这相对容易(Buntine 1996),并且在图模型中,可以训练条件概率表或它们的参数设置(如式(16.30)中的参数),以最大化似然,或者如果已知适当的先验,可以用贝叶斯方法(参见第 14 章)。

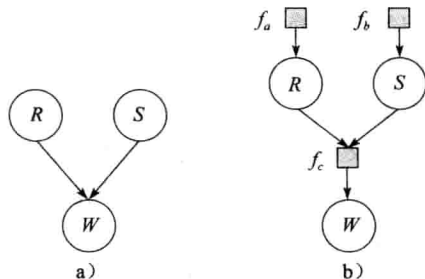


图 16-15 a) 一个有向图,教化后它将有一个环;b)它对应的因素图是一棵结树。三个因素是 $f_a(R) \equiv P(R)$, $f_b(S) \equiv P(S)$, $f_c(R, S, W) \equiv P(W|R, S)$

第二部分更困难、更有趣,是学习图结构(Cowell 等 1999)。这基本上是一个模型选择问题,就像学习多层感知器结构的增量方法一样(参见 11.9 节),我们可以将这部分看作在所有可能的图空间中进行搜索。例如,我们可以考虑增加或删除一条弧、增加或删除一个隐藏结点的操作,然后进行搜索、(在每次中间迭代使用参数学习)评估每一步的改进。然而需要注意,为了检查过拟合,我们应该适当地正则化,这对应于偏爱较简单图的贝叶斯方法(Neapolitan 2004)。然而,由于状态空间很大,如果有人类专家,能够手工地定义变量之间的因果关系,并创建变量小的组群子图,将是特别有益的。

16.8 影响图

正如在第 3 章中我们将概率推广到具有风险的动作一样,影响图(influence diagram)是一种图形模型,它是包括决策和效用的图模型的推广。一个影响图包含机会节点(chance node),表示我们在图模型中使用的随机变量(参见图 16-16)。影响图还包含决策节点和效用节点。决策节点(decision node)表示动作的选择。效用节点(utility node)用于计算效用。决策可以根据机会节点做出,并且可能影响其他机会节点和效用节点。

影响图上的推断是对图模型上的信念传播的扩展。给定一些机会节点上的证据,证据传播,并且对每一个可能的决策计算效用,选择具有最大效用的决策,对一个给定输入分类的影响图如图 16-16 所示。给定输入,决策节点决定类,对于每个选择,我们得到一定的效用(风险)。

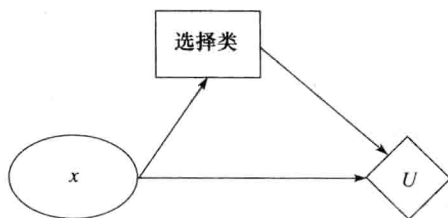


图 16-16 对应于分类的影响图。依赖于输入 x , 选择一个导致一定效用(风险)的类

16.9 注释

图模型有两个优点。一个优点是我们可以呈现变量的相互影响,更好地理解过程,例如,通过使用因果生成模型。第二个优点是通过找出对应于诸如贝叶斯规则或边缘化等基本概率过程的图操作,可以将推断任务映射到可以有效表示和实现的通用图算法。

变量和它们之间的依赖性用图直观地表示,以及许多变量的复杂全局函数分解成每个都只涉及变量的一个小子集的局部函数的乘积的思想,被用在决策、编码和信号处理的不同领域,Kschischang、Frey 和 Loeliger(2001)给出了评述。

推断算法在多树或结树上的复杂度由父节点的最大个数或最大团的大小决定。当这些很大时,严格推断可能是不可行的。在这种情况下,我们需要使用近似或抽样算法。Jordan 等 1999、MacKay 2003、Andieu 等 2003、Bishop 2006 和 Jordan 2009 讨论了各种近似算法和马尔可夫链蒙特卡罗(Markov chain Monte Carlo, MCMC)算法。

图模型特别适合表示贝叶斯方法,其中除了代表变量节点之外,还有代表影响观测变量的隐藏参数的节点。我们还可以引进分层结构,其中有代表超参数的节点(即代表一级参数先验的二级参数等)。

隐马尔可夫模型也是一种图模型,并且事实上任何图模型都可以在时间上展开并加上相继复制之间的依赖性,从而在时间上扩展。这种动态图模型(dynamic graphical model)在也包

括时间维的领域找到了应用——例如，语音识别。事实上，隐马尔可夫模型只不过是一种聚类序列问题，其中时刻 t 的簇索引不仅依赖于时刻 t 的观测，而且依赖于时刻 $t-1$ 的索引，并且 Baum-Welch 算法是期望最大化算法的扩展，以便也包括这种时间上的依赖性。在 6.4 节中我们讨论过因素分析，其中少量的隐藏因素生成了观测；类似地，线性动态系统可以看作这种因素分析模型的序列，其中当前因素也依赖以前的因素。

需要时可以加上动态依赖性。例如，图 16-15 对特定的一天湿草地的原因建模；如果我们相信昨天的天气对今天的天气有影响（并且我们应该相信——在相继的几天趋向于多云，然后几天晴天，如此下去），则我们可以有图 16-17 所示的动态图模型，其中我们对这种依赖性建模。

一般的图模型机制使得我们超越纯粹的 HMM 能力并且导致性能的提高，例如，在语音识别中 (Zweig 2003, Bilmes 和 Bartels 2005)。图模型也可以用于计算机视觉，例如，信息检索 (Barnard 等 2003) 和场景分析 (Sudderth 等 2008)。图模型在生物信息学中的应用的一个综述在 Donkers 和 Tuyls 2008 中。

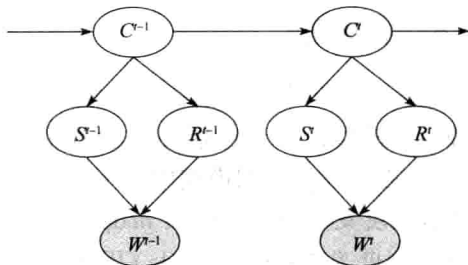


图 16-17 一个动态版本，其中有一个展示相继几天天气依赖性的图的链

16.10 习题

1. 在分类问题中使用两个独立的输入变量，即 $p(x_1, x_2 | C) = p(x_1 | C)p(x_2 | C)$ ，如何计算 $p(x_1 | x_2)$ ？为 $p(x_j | C_i) \sim \mathcal{N}(\mu_{ij}, \sigma_{ij}^2)$ 推导公式。
2. 对于头到头节点，证明式 (16.10) 蕴含 $P(X, Y) = P(X) \cdot P(Y)$ 。
3. 在图 16-4 中，计算 $P(R | W)$ 、 $P(R | W, S)$ 和 $P(R | W, \sim S)$ 。
4. 在式 (16.30) 中， X 是二元的。如果 X 可以取 K 个离散值之一，我们需要对它做何修改？
5. 证明：在联合分布可以表示成式 (16.12) 的有向图中， $\sum_{\mathbf{x}} p(\mathbf{x}) = 1$ 。
6. 将 Necker 立方体绘制成定义链以指示互斥地增强或禁止不同拐角解释间的关系。
7. 如何在图 16-17 所示的动态天气图上进行推断？
8. 用图 16-9 的方式，为两类的线性逻辑斯蒂回归给出图模型。
9. 提出一种合适的度量，这种度量可以在学习图结构中用作状态搜索。合适的操作是什么？
10. 一般而言，在报纸上，一个报道者在连续几天内写一系列关于同一主题的文章作为新闻报道追踪。如何使用图模型对其建模？

16.11 参考文献

- Andrieu, C., N. de Freitas, A. Doucet, and M. I. Jordan. 2003. "An Introduction to MCMC for Machine Learning." *Machine Learning* 50: 5-43.
- Barnard, K., P. Duygulu, D. Forsyth, N. de Freitas, D. M. Blei, and M. I. Jordan. 2003. "Matching Words and Pictures." *Journal of Machine Learning Research* 3: 1107-1135.
- Bilmes, J., and C. Bartels. 2005. "Graphical Model Architectures for Speech Recognition." *IEEE Signal Processing Magazine* 22: 89-100.

- Bishop, C. M. 2006. *Pattern Recognition and Machine Learning*. New York: Springer.
- Buntine, W. 1996. "A Guide to the Literature on Learning Probabilistic Networks from Data." *IEEE Transactions on Knowledge and Data Engineering* 8: 195-210.
- Cowell, R. G., A. P. Dawid, S. L. Lauritzen, and D. J. Spiegelhalter. 1999. *Probabilistic Networks and Expert Systems*. New York: Springer.
- Donkers, J., and K. Tuyls. 2008. "Belief Networks in Bioinformatics." In *Computational Intelligence in Bioinformatics*, ed. A. Kelemen, A. Abraham, and Y. Chen, 75-111. Berlin: Springer.
- Ghahramani, Z. 2001. "An Introduction to Hidden Markov Models and Bayesian Networks." *International Journal of Pattern Recognition and Artificial Intelligence* 15: 9-42.
- Jensen, F. 1996. *An Introduction to Bayesian Networks*. New York: Springer.
- Jordan, M. I., ed. 1999. *Learning in Graphical Models*. Cambridge, MA: MIT Press.
- Jordan, M. I. 2004. "Graphical Models." *Statistical Science* 19: 140-155.
- Jordan, M. I. 2009. *An Introduction to Probabilistic Graphical Models*. Forthcoming.
- Jordan, M. I., Z. Ghahramani, T. S. Jaakkola, and L. K. Saul. 1999. "An Introduction to Variational Methods for Graphical Models." In *Learning in Graphical Models*, ed. M. I. Jordan, 105-161. Cambridge, MA: MIT Press.
- Kschischang, F. R., B. J. Frey, and H.-A. Loeliger. 2001. "Factor Graphs and the Sum-Product Algorithm." *IEEE Transactions on Information Theory* 47: 498-519.
- Lauritzen, S. L., and D. J. Spiegelhalter. 1988. "Local Computations with Probabilities on Graphical Structures and their Application to Expert Systems." *Journal of Royal Statistical Society B* 50: 157-224.
- MacKay, D. J. C. 2003. *Information Theory, Inference, and Learning Algorithms*. Cambridge, UK: Cambridge University Press.
- Neapolitan, R. E. 2004. *Learning Bayesian Networks*. Upper Saddle River, NJ: Pearson.
- Pearl, J. 1988. *Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference*. San Francisco, CA: Morgan Kaufmann.
- Pearl, J. 2000. *Causality: Models, Reasoning, and Inference*. Cambridge, UK: Cambridge University Press.
- Sudderth, E. B., A. Torralba, W. T. Freeman, and A. S. Willsky. 2008. "Describing Visual Scenes Using Transformed Objects and Parts." *International Journal of Computer Vision* 77: 291-330.
- Zweig, G. 2003. "Bayesian Network Structures and Inference Techniques for Automatic Speech Recognition." *Computer Speech and Language* 17: 173-193.

组合多学习器

我们在前面的章节中讨论了许多不同的学习算法。尽管一般而言它们是成功的，但没有哪一个算法总是最准确的。现在，我们将讨论由多个学习器组成的模型。这些学习器互补，因此通过组合它们，我们可以获得更高的准确率。

17.1 基本原理

在任何应用中，我们可以使用多个学习算法中的一个，而使用某算法，存在对最终学习器产生影响的超参数。例如，在分类情况下，我们可以使用参数分类器或多层感知器，而比如对一个多层感知器，我们还要确定隐藏单元的数目。“没有免费的午餐”法则表明没有一个学习算法可以在任何领域总是产生最准确的学习器。通常的方法是试验很多种算法，然后选择一个在单独的验证集上性能最佳的算法。

每一个学习算法都构建了一个基于一组假设的某种模型。当假设在数据上不成立时，这种归纳偏倚将导致误差。学习是一个不定问题，并且在有限的数据上，每个学习算法都收敛到不同的解，并在不同的情况下失效。可以通过性能调节使得一个学习算法在验证集上达到尽可能最高的准确率，但是调节本身就是一个复杂的任务，并且即使对最好的学习器而言，也存在一些实例使其无法足够准确。解决之道在于也许存在另一种学习方法，在这些实例上是准确的。通过合适方式将多个学习器组合可以提高准确率。近来，随着计算和存储变的更为廉价，组合多个学习器的系统也随之流行。

这里有两个基本问题：

- 1) 如何产生互补的基学习器？
- 2) 为了最大化准确率，如何组合基学习器的输出？

本章的讨论将回答这两个相关的问题。我们将会看到，模型组合并非总是能够提高准确率的诀窍；模型组合确实总是增加训练和检验的时间和空间复杂度，并且除非小心地训练基学习器并且精明地组合它们的决策，否则我们只会为这种附加的复杂性付出代价，而在准确率方面得不到显著收益。

17.2 产生有差异的学习器

由于组合总是采取类似决策的学习器是没有意义的，因此我们的目标是可以寻得一组有差异的 (diverse) 学习器，它们采取不同的决策，使得相互之间可以互补。与此同时，除非这些学习器至少在它们的专门领域是准确的，否则组合也不可能获得整体成功。因此，我们具有双重任务：最大化个体学习器的准确率和学习器之间的差异性。现在，我们讨论实现这些目标的不同方法。

1. 不同的算法

我们可以使用不同的学习算法来训练不同的基学习器。不同的算法对数据做不同的假设并导致不同的分类器。例如，一个基学习器可以是参数的，而另一个可以是非参数的。当我们决定使用一个算法的时候，是将重点放在单一方法上并忽略所有其他方法。通过组合基于多个算法的多个学习器，我们摆脱了做单一决策，并且再也不将所有鸡蛋放在同一个篮子中。

2. 不同的超参数

我们可以使用相同的学习算法，但使用不同的超参数。这样的例子包括：多层感知器中的隐藏单元数目， k -最近邻中的 k 值，决策树中的误差阈值，支持向量机中的核函数等。对高斯参数分类器而言，协方差矩阵是否共享是一个超参数。如果在优化算法中使用诸如梯度下降这样的最终状态依赖于初始状态的迭代过程，如使用向后传播的多层感知器，则初始状态(如初始权重)是另一种超参数。当我们用不同的超参数值训练多个基学习器时，我们对这种因素取平均值并降低方差，从而减小误差。

3. 不同的输入表示

不同的基学习器也可能使用相同输入对象或事件的不同表示，从而使得集成不同类型的感知器、测量、模态成为可能。不同的表示凸显了对象的不同特征，从而产生更好的识别。在许多应用中，存在多个信息源，我们期望使用所有这些数据来提取更多的信息，并在预测中到达更高的准确率。

例如，在语音识别中，为了识别语音单词，除了声学输入之外，我们还可以使用讲话者在读出单词时的嘴唇视频图像。这类似于传感器融合(sensor fusion)，其中来自不同传感器的数据集成在一起，为特定应用提取更多的信息。

最简单的方法就是连接所有数据向量并将其当作是来自同一数据源的一个大向量，但是这种方法在理论上似乎不太合适，因为这样相当于对取样于元统计分布的数据进行建模。此外，更高的输入维度使得系统更为复杂，并且为了得到精确的估计需要更大的样本。我们采用的方法是使用不同的基学习器在不同的数据源上分别进行预测，然后组合它们的预测结果。

即使只有一种输入表示，即通过从中选择随机子集，我们也可以有使用不同输入特征的多个分类器，这称为随机子空间方法(random subspace method) (Ho 1998)。这具有如下效果：不同的学习器将从不同的视角考察同一问题，并且将是鲁棒的；这也有助于降低维灾难，因为输入只有不多的维。

4. 不同的训练集

另一种可能的方法是使用训练集的不同子集来训练不同的基学习器，这可以通过在给定的样本上随机抽取随机训练集来实现，这称作装袋(bagging)。抑或，可以串行地训练学习器，使得前一个基学习器上预测不准的实例在之后的基学习器训练中获得更多的重视，这种

419
420

例子有提升(boosting)和级联(cascading), 这些方法积极地尝试生成互补的学习器, 而不是靠碰运气。

训练样本的划分也可以基于输入空间的局部性来完成, 使得每个基学习器在输入空间的某一局部的实例上训练, 这就是我们在第 12 章讨论的“混合专家模型”所做的, 不过我们从多学习器组合的角度来再次讨论。类似地, 可以将主任务定义为由基学习器实现的若干子任务, 如纠错输出码(error-correcting output code)所做的那样。

5. 差异性与准确率

非常重要的一点是当生成多个基学习器时, 只要它们有合理的准确率即可, 而不要求它们每个都非常准确, 因此不要, 也不必对这些基学习器进行单独优化以获取最佳准确率。基学习器的选择并不是由于它的准确性, 而是由于它的简单性。然而, 我们的确要求基学习器是有差异的, 即在不同实例上是准确的, 专注于问题的子领域。我们所关心的是基学习器组合后的准确性, 而非开始时各基学习器的准确性。例如, 我们有一个准确率为 80% 的分类器。当确定第二个分类器的时候, 我们不关心其总体准确率, 只要我们知道何时使用哪个分类器, 就可以只关心其在第一个分类器误分类的 20% 实例上的准确性如何。

正如我们将要讨论的, 这意味着所要求的学习器的准确性和差异性也依赖于如何组合它们的决策。如果像投票策略中那样, 学习器用于所有输入, 则它应该处处准确, 并且必须处处存在差异。如果我们将输入空间划分成针对不同学习器的专门区域, 则差异性就已经被划分所保证, 而学习器只需要在它的局部域中是准确的即可。

17.3 模型组合方案

存在不同的组合多个基学习器来产生最终输出的方法:

- 多专家组合(multiexpert combination)方法让基学习器并行工作。这些方法可以进一步划分成两类:
 - 全局(global)方法又称学习器融合(learner fusion)。在全局方法中, 给定一个输入, 所有的基学习器都产生一个输出, 并且所有这些输出都要使用, 例子包括投票(voting)和层叠(stackings)。
 - 在局部(local)方法, 或学习器选择(learner selection)中, 例如在混合专家(mixture of experts)中, 有一个选通(gating)模型, 它考察输入, 并选择一个(或几个)学习器来产生输出。
- 多级组合(multistage combination)方法使用一种顺序(serial)方法, 其中下一个基学习器只在前一个基学习器预测不够准确的实例上进行训练或检验。其基本思想是, 基学习器(或其所使用的不同表示)按复杂度递增排序, 使得除非前一个更简单的基学习器的结果是不足够可信的, 否则就不使用复杂的基学习器(或不提取其复杂表示), 一个这样的例子是级联(cascading)。

假设有 L 个基学习器。我们用 $d_j(x)$ 表示基学习器 \mathcal{M}_j 在给定的任意维输入 x 上的预测。在存在多种输入数据表示的情况下, 每个 \mathcal{M}_j 使用一种不同的输入表示 x_j 。最后的预测由各个基学习器的预测计算:

$$y = f(d_1, d_2, \dots, d_L | \Phi)$$
 (17.1)

其中 $f(\cdot)$ 是一个组合函数， Φ 表示其参数。

当有 K 个输出的时候，每个学习器有 K 个输出 $d_{ji}(x)$ ， $i=1, \dots, K, j=1, \dots, L$ ，而组合它们，我们仍然产生 K 个值 $y_i, i=1, \dots, K$ 。进而，比如在分类中，我们选择具有最大 y_i 值的类：

如果 $y_i = \max_{k=1}^K y_k$ ，则选择 C_i 。

17.4 投票法

组合多个分类器的最简单方法是通过投票(voting)，这相当于取学习器的线性组合(参见图 17-1)：

$$y_i = \sum_j w_j d_{ji}, \quad \text{其中 } w_j \geq 0, \sum_j w_j = 1$$
 (17.2)

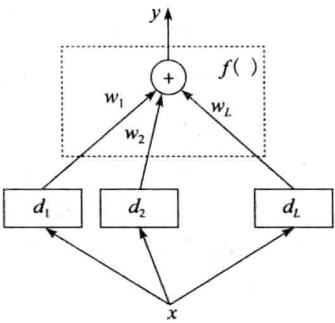


图 17-1 基学习器是 d_j ，而它们的输出用 $f(\cdot)$ 组合。这是单输出；对于分类，每个基学习器都有 K 个输出分别用于计算 y_j ，然后我们选择最大的。注意，图中所有学习器观测相同的输入；可能存在不同的学习器观测相同输入对象或事件的不同表示

这种方法也称系综(ensemble)或线性集思广益(linear opinion pool)。在最简单的情况下，所有的学习器都赋予相等的权重，而我们有简单投票(simple voting)，对应于取平均值。取(加权)和也只是一种可能的方法，还有一些其他组合规则，如表 17-1 所示(Kittler 等 1998)。如果输出不是后验概率，则这些规则要求输出规范化为相同的尺度(Jain, Nandakumar 和 Ross 2005)。

423
424

表 17-1 分类器组合规则

规则	融合函数 $f(\cdot)$
和	$y_i = \frac{1}{L} \sum_{j=1}^L d_{ji}$
加权和	$y_i = \sum_j w_j d_{ji}, w_j \geq 0, \sum_j w_j = 1$
中位数	$y_i = \text{median}_j d_{ji}$
最小值	$y_i = \min_j d_{ji}$
最大值	$y_i = \max_j d_{ji}$
乘积	$y_i = \prod_j d_{ji}$

使用这些规则的一个例子显示在表 17-2 中, 它展示了不同规则的效果。求和规则最直观, 并且实践中使用最广泛。中位数规则对离群点更鲁棒; 最小和最大规则分别是悲观和乐观的。使用乘积规则, 每个学习器都有否决权; 无论其他学习器如何投票, 如果一个学习器的输出为 0, 则整体输出就为 0。注意, 使用组合规则之后, y_i 之和不必为 1。

表 17-2 在三个学习器和三个类上组合规则的例子

规则	C_1	C_2	C_3
d_1	0.2	0.5	0.3
d_2	0.0	0.6	0.4
d_3	0.4	0.4	0.2
和	0.2	0.5	0.3
中位数	0.2	0.5	0.4
最小值	0.0	0.4	0.2
最大值	0.4	0.6	0.4
乘积	0.0	0.12	0.032

在加权和中, d_{ji} 是学习器 j 对 C_i 类的投票, w_j 是它的投票权重。简单投票是加权投票的特殊情况, 其中所有的投票者具有相同的权重, 即 $w_j = 1/L$ 。在分类中, 这称为简单多数表决 (plurality voting), 其中得票最多的类胜出。当只有两个类时, 这就是多数表决 (majority voting), 其中胜出类获取一半以上的票 (习题 1)。如果投票者还能提供它们为每个类投票多少的额外信息 (例如, 通过后验概率), 则规范化后, 这些信息即可用作加权投票 (weighted voting) 方案的权重。同样, 如果 d_{ji} 是类后验概率 $P(C_i | x, \mathcal{M}_j)$, 则我们可以简单将其相加 ($w_j = 1/L$) 并选取具有最大 y_i 的类。

对于回归, 可以使用简单平均、加权平均或中位数来融合基回归器的输出。中位数对噪声比平均值更鲁棒。

另外一种找出 w_j 的可能方法是在另外的验证集上评估学习器 (回归器或分类器) 的准确率并使用这些信息来计算权重, 使得我们可以对更准确的学习器赋予更高的权重。这些权重也可以由数据来学习, 我们将在 17.9 节讨论层叠泛化时讨论。

投票方案可以看作贝叶斯框架下的近似, 以权重近似先验模型概率, 并以模型决策近似模型条件似然, 这就是贝叶斯模型组合 (Bayesian model combination)。例如, 在分类中我们有 $w_j = P(\mathcal{M}_j)$, $d_{ji} = P(C_i | x, \mathcal{M}_j)$, 而式 (17.2) 对应于:

$$P(C_i | x) = \sum_{\text{所有模型 } \mathcal{M}_j} P(C_i | x, \mathcal{M}_j) P(\mathcal{M}_j) \quad (17.3)$$

简单投票相当于假定一致先验概率。如果先验分布更倾向于较简单的模型, 则简单投票将赋予简单模型更大的权重。我们可以不集成所有模型, 只选取一个我们认为 $P(\mathcal{M}_j)$ 值高的子集, 或者可以执行另一个贝叶斯步骤并计算给定样本上模型的条件概率 $P(\mathcal{M}_j | \mathcal{X})$, 并从该密度中选取一些高概率的模型。

Hansen 和 Salamon (1990) 证明: 给定 (预测) 成功概率高于 1/2 (即比随机猜测的预测好) 的一组独立的两类分类器, 使用多数表决, 预测准确率随着投票分类器个数的增加而提高。

假设 d_j 是独立同分布的, 其期望值为 $E[d_j]$ 、方差为 $\text{Var}(d_j)$, 那么当使用 $w_j = 1/L$ 取简单平均时, 输出的期望值和方差分别为:

$$E[y] = E\left[\sum_j \frac{1}{L}d_j\right] = \frac{1}{L}LE[d_j] = E[d_j]$$

$$\text{Var}(y) = \text{Var}\left(\sum_j \frac{1}{L}d_j\right) = \frac{1}{L^2}\text{Var}\left(\sum_j d_j\right) = \frac{1}{L^2}L\text{Var}(d_j) = \frac{1}{L}\text{Var}(d_j) \quad (17.4)$$

可以看到期望值没有改变,因而偏倚也不改变。但是方差,还包括均方误差随着独立投票者数量 L 的增加而降低。在一般情况下,

$$\text{Var}(y) = \frac{1}{L^2}\text{Var}\left(\sum_j d_j\right) = \frac{1}{L^2}\left(\sum_j \text{Var}(d_j) + 2\sum_j \sum_{i < j} \text{Cov}(d_j, d_i)\right) \quad (17.5)$$

这暗示,如果学习器是正相关的,则方差(和误差)增加。这样,我们可以将使用不同的算法和输入特征看作,如果不完全消除正相关性,也是在努力减少正相关性。在 17.10 节,我们将讨论从系综中删除高度正相关的学习器的剪枝方法。

这里,我们还看到,如果投票者并非独立而是负相关的,则进一步降低方差是可能的。如果随此增加的偏倚不是更高的话(因为这些目标是矛盾的),则误差也会降低。我们不可能有大量分类器,它们都是准确的并且是负相关的。例如,在混合专家模型中,学习器是局部化的,专家是负相关的,但是有偏倚的(Jacobs 1997)。

如果将每个基学习器看作是附加在真实判别式或回归函数上的随机噪声函数,而且这些噪声函数是不相关的并且均值为 0,那么在每个估计上的平均就类似于噪声上的平均。从这种意义上讲,投票具有光滑函数空间的效果,并且可以看作是一个在真实函数上具有光滑假设的正则化子(Perrone 1993)。我们在图 4-5(d)中看到一个例子,其中通过在具有大方差的模型上取平均,我们得到了比单独模型更好的拟合。这就是投票的思想:我们对具有高方差低偏倚的模型投票,使得在组合后,偏倚依然保持很小而通过取平均则降低了方差。即使单个模型是有偏倚的,方差的降低也可能抵消偏倚并且仍然可能降低误差。

17.5 纠错输出码

在纠错输出码(Error-Correcting Output Code, ECOC)中(Dietterich 和 Bakiri 1995),主要的分类任务通过由基学习器实现的一组子任务来定义。其思想是:将一个类从其他类区分开来的原始任务可能是一个困难的问题。作为替代,我们可以定义一组简单的分类问题,每个专注于原始任务的一个方面,并通过组合这些简单分类器来得到最终的分类器。

这时,基学习器是输出为 $-1/+1$ 的二元分类器,并且有一个 $K \times L$ 的编码矩阵 \mathbf{W} ,其 K 行是关于 L 个基学习器 d_j 的类的二元编码。例如,如果 \mathbf{W} 的第二行是 $[-1, +1, +1, -1]$,则这意味如果一个实例属于 C_2 类,则该实例应在 d_1 和 d_4 上取负值,在 d_2 和 d_3 上取正值。类似地,编码矩阵的列定义了基学习器的任务。例如,如果第三列是 $[-1, +1, +1]^T$,可理解为第三个基学习器 d_3 的任务是将属于 C_1 类的实例与属于 C_2 和 C_3 类的实例区分开。这就是我们如何构成基学习器的训练集的方式。例如,在这个例子中,所有标记为 C_2 类及 C_3 类的实例形成 X_3^+ ,而标记为 C_1 类的实例构成 X_3^- ,而对 d_3 的训练应使得当 $x' \in X_3^+$ 时输出 $+1$,当 $x' \in X_3^-$ 时输出 -1 。

357

这样,编码矩阵使得我们可以用二分问题($K=2$ 的分类问题)定义多分问题($K>2$ 的分类问题),并且这是一种适用于任意可以实现二分基学习器的学习算法的方法,例如,线性或多层感知器(单输出)、决策树或初始定义用于两类问题的 SVM。

典型的每类一个判别式的情况对应于对角编码矩阵, 其中 $L = K$ 。例如, 对于 $K = 4$, 我们有

$$W = \begin{bmatrix} +1 & -1 & -1 & -1 \\ -1 & +1 & -1 & -1 \\ -1 & -1 & +1 & -1 \\ -1 & -1 & -1 & +1 \end{bmatrix}$$

这里的问题是: 如果某一个基学习器存在错误, 就会有误分类, 因为类的码字之间非常相似。因而纠错码采用的方法是使 $L > K$ 来增加码字之间的汉明距离。一种可能的方法是类逐对分开 (pairwise separation), 其中对 $i < j$ 有一个不同的基学习器将 C_i 与 C_j 分开 (见 10.4 节)。在这种情况下, 当 $K = 4$ 时 $L = K(K-1)/2$, 编码矩阵为

$$W = \begin{bmatrix} +1 & +1 & +1 & 0 & 0 & 0 \\ -1 & 0 & 0 & +1 & +1 & 0 \\ 0 & -1 & 0 & -1 & 0 & +1 \\ 0 & 0 & -1 & 0 & -1 & -1 \end{bmatrix}$$

其中的 0 表示“无关”。这就是说, 训练 d_1 来将 C_1 与 C_2 分开并且在训练中不使用属于其他类的实例。类似地, 一个实例属于 C_2 如果有 $d_1 = -1$ 并且 $d_4 = d_5 = +1$, 并且我们不考虑 d_2 、 d_3 和 d_6 的值。这种方法的问题是 L 是 $O(K^2)$ 。因而, 对于比较大的 K , 逐对分开可能是不可行的。

方法是预先设置 L 值, 然后寻找 W 使得以汉明距离衡量的行间距离以及列间距离都尽可能的大。对 K 个类而言, 存在 $2^{(K-1)} - 1$ 种可能的列, 即两类问题。这是因为 K 位可写为 2^K 种不同的形式和补 (比如, “0101” 和 “1010”, 从我们的角度来看, 二者定义相同的判别式), 将有可能组合除以 2 然后减 1, 因为全为 0 (或 1) 的列是无用的。例如, 当 $K = 4$ 时, 我们有

$$W = \begin{bmatrix} -1 & -1 & -1 & -1 & -1 & -1 & -1 \\ -1 & -1 & -1 & +1 & +1 & +1 & +1 \\ -1 & +1 & +1 & -1 & -1 & +1 & +1 \\ +1 & -1 & +1 & -1 & +1 & -1 & +1 \end{bmatrix}$$

当 K 很大时, 对于一个给定的 L 值, 我们从 $2^{(K-1)} - 1$ 列选取 L 列。我们希望 W 的这些列尽可能的不相同, 以便每个基学习器所学习的子任务尽可能互不相同。同时, 我们希望 W 的行也尽可能不相同, 使得在一个或多个基学习器失效时可以获得最大的纠错。

ECOC 可以投票方式来表述, 其中 W 的元素 w_{ij} 被看作是投票权重:

$$y_i = \sum_{j=1}^L w_{ij} d_j \quad (17.6)$$

然后我们选取具有最高 y_i 的类。通过求加权和并选择最大值 (判断类别) 取代寻求一个精确的匹配使得 d_j 也不必是二元的, 而是可取 -1 到 $+1$ 之间的任意值, 以软确定性取代硬判决。注意, 位于 0 和 1 之间的 p_j 值 (例如后验概率) 可以很简单地被转换为 -1 到 $+1$ 间

的 d_j 值:

$$d_j = 2p_j - 1$$

式(17.6)与式(17.2)的一般投票模型的不同在于投票的权重对不同的类可以不同,即以 w_{ij} 取代了 w_j , 并且 $w_j \geq 0$ 而 w_{ij} 为 -1 、 0 或 $+1$ 。

ECOC 的一个问题是: 由于编码矩阵 \mathbf{W} 被设置为先验, 因此不能保证由 \mathbf{W} 的列所定义的任务一定是简单的。Dietterich 和 Bakiri (1995) 的研究表明二分树可能要比多分树大, 而且当使用多层感知器时, 后向传播可能收敛较慢。

17.6 装袋

装袋 (bagging) 是一种投票方法, 其中基学习器通过在稍有差异的训练集上训练而有所不同。从给定的样本数据上产生 L 个稍稍不同的样本集通过自助法来完成, 其中给定一个大小为 N 的数据集 \mathcal{X} , 随机从 \mathcal{X} 中有放回地 (with replacement) 抽取 N 个实例。由于抽样是通过有放回方式完成的, 可能某些实例被多次抽取而某些实例根本没有被抽到。当抽取 L 个样本 $\mathcal{X}_j (j=1, \dots, L)$ 后, 这些样本集是彼此相似的, 因为它们是从相同的原始样本数据源抽取的, 但是源于随机性而又稍有不同。基学习器 d_j 在这 L 个样本集的 \mathcal{X}_j 上训练。

一个学习算法是不稳定算法 (unstable algorithm), 如果训练集中很小的变化会引起所产生学习器很大的差异, 即学习算法具有高方差。装袋是自助聚集 (bootstrap aggregating) 的简单说法, 就是使用自助法产生 L 个训练集, 并使用不稳定的学习过程训练 L 个基学习器, 并在检验时取 (预测的) 平均值 (Breiman 1996)。装袋可用于分类和回归。在用于回归的情况下, 为了更加鲁棒, 可以在组合预测结果时以中值来取代平均值。

前面, 我们已经看到, 仅当正相关性小时, 取平均值才能降低方差; 如果算法在相同数据集的再抽样版本上多次运行导致具有高正相关性的学习器, 则该算法是稳定的。

诸如决策树和多层感知器这样的算法是不稳定的。最近邻算法是稳定的, 但是精简的最近邻算法是不稳定的 (Alpaydm 1997)。如果原始训练集很大, 则我们可能希望使用自助法来从它产生小一些的数据集 ($N' < N$), 否则 \mathcal{X}_j 的自助副本将会非常相似, 从而 d_j 将高度相关。

17.7 提升

在装袋中, 产生互补的基学习器是靠运气以及学习方法的不稳定性。在提升中, 我们通在前一个学习器所犯的误差上训练下一个学习器, 积极地尝试产生互补的学习器。原始的提升 (boosting) 算法 (Schapire 1990) 组合了三个弱学习器来产生一个强学习器。所谓弱学习器 (weak learner) 是误差概率小于 $1/2$ 的学习器, 这使得它对两类问题比随机猜测要好, 而强学习器 (strong learner) 具有任意小的误差概率。

给定一个大训练集, 我们随机地将其划分为三部分。使用 \mathcal{X}_1 来训练 d_1 。然后取 \mathcal{X}_2 并将它馈入 d_1 。将所有 d_1 误分类的实例以及在 \mathcal{X}_2 中被 d_1 正确分类的一些实例一起作为 d_2 的训练集。然后取 \mathcal{X}_3 并将它馈入 d_1 和 d_2 。其中 d_1 和 d_2 输出不一致的实例形成 d_3 的训练集。在检验期间, 给定一个实例我们首先将其提供给 d_1 和 d_2 ; 如果二者输出一致, 这就作为输出

结果, 否则 d_3 的输出作为结果。Schapire(1990)的研究表明这个整体系统降低了错误率, 并且错误率可以通过递归地使用这样的系统(即将三个模型构成的提升系统作为更高层系统的 d_j)而任意降低。

尽管这种方法很成功, 但是提升方法的不足之处在于需要一个非常大的训练集。样本需要一分为三, 而且第二和第三个分类器只在由其前的分类器犯错的实例构成的子集上训练。因此, 除非有一个很大的训练集, 否则 d_2 和 d_3 将无法拥有合理大小的训练集。Drucker 等(1994)在其提出的提升多层感知器中使用了有 118 000 个实例的数据集, 用于光学手写数字识别。

431 Freund 和 Schapire(1996)提出了提升的一个变种, 叫 AdaBoost, 是自适应提升的缩写, 其中重复使用相同的训练集因而而不要求数据集很大。AdaBoost 还可以组合任意数量的基学习器, 不一定是三个。

已经有很多 AdaBoost 的变种被提出; 这里我们讨论原始的算法 AdaBoost.M1(见图 17-2): 其思想是将实例抽取的概率修改成误差的函数。令 p_j^i 表示实例对 (x^i, r^i) 被抽取用于训练第 j 个基学习器的概率。最初, 所有的 $p_1^i = 1/N$ 。然后, 以如下方式添加新的基学习器: 从 $j=1$ 开始, ε_j 表示 d_j 的错误率。AdaBoost 要求任意的 $\varepsilon_j < 1/2$; 如果不满足, 即停止添加新的基学习器。注意, 这里的错误率并非基于原始问题, 而是基于在第 j 步中使用的数据集。定义 $\beta_j = \varepsilon_j / (1 - \varepsilon_j) < 1$, 并且设置 $p_{j+1}^i = \beta_j p_j^i$ 如果 d_j 正确地对 x^i 分类, 否则设置 $p_{j+1}^i = p_j^i$ 。由于 p_{j+1}^i 应该是概率, 所以我们用 p_{j+1}^i 除以 $\sum_i p_{j+1}^i$ 对其规范化, 使它们的和为 1。这样做的效果是将被正确分类的实例的(抽取)概率降低, 而将被错误分类的实例的概率提高。然后, 根据这些修改后的概率 p_{j+1}^i , 从原样本中有放回地抽取相同大小的样本集, 并用于训练 d_{j+1} 。

这样做的效果是使得 d_{j+1} 更专注于被 d_j 误分类的实例。这就是为什么基学习器以简单而不是准确为原则选取, 否则下一个训练样本集将仅仅包含少数离群点和噪声实例的多次重复。例如, 对于决策树, 使用的是决策树桩(decision stump), 一种只有一层或两层的树。因此, 很明显它们是有偏的但是方差上的降低比较大, 而且总体误差也会降低。像线性判别式这样的算法本身具有低方差, 我们不能通过 AdaBoost 而获得更低的方差。

一旦完成训练, AdaBoost 就采用投票方法。给定一个实例, 所有的 d_j 决定其分类, 而后取一个加权的投票结果, 其中权重与基学习器(在训练集上的)准确率成正比: $w_j = \log(1/\beta_j)$ 。Freund 和 Schapire(1996)表明 AdaBoost 在 22 个基准问题上提高了准确率, 在一个基准问题上准确率相同, 而在 4 个基准问题上准确率较差。

训练:

For 所有的 $\{x^i, r^i\}_{i=1}^N \in \mathcal{X}$, 初始化 $p_1^i = 1/N$

For 所有的基分类器 $j=1, \dots, L$

按照概率 p_j^i 随机地从 \mathcal{X} 抽取 x_j

使用 x_j 训练 d_j

For 每个 (x^i, r^i) , 计算 $y_j^i \leftarrow d_j(x^i)$

计算错误率: $\varepsilon_j \leftarrow \sum_i p_j^i \cdot 1(y_j^i \neq r^i)$

If $\varepsilon_j > 1/2$, then $L \leftarrow j-1$; Stop

$\beta_j \leftarrow \varepsilon_j / (1 - \varepsilon_j)$

For 每个 (x^i, r^i) // 如果正确, 则减低概率

If $y_j^i = r^i$ then $p_{j+1}^i \leftarrow \beta_j p_j^i$ else $p_{j+1}^i \leftarrow p_j^i$

// 规范化概率:

$Z_j \leftarrow \sum_i p_{j+1}^i$; $p_{j+1}^i \leftarrow p_{j+1}^i / Z_j$

检验:

给定 x , 计算 $d_j(x)$, $j=1, \dots, L$

计算类输出, $i=1, \dots, K$:

$$y_i = \sum_{j=1}^L \left(\log \frac{1}{\beta_j} \right) d_{ji}(x)$$

图 17-2 AdaBoost 算法

Schapire 等(1998)认为 AdaBoost 的成功源于其扩展边缘(margin)。如果边缘增加,训练实例可以更好的被分隔而使误分类不易发生。这使得 AdaBoost 的目标和支撑向量机(参见 13 章)类似。

在 AdaBoost 中,尽管不同的基学习器使用稍有差异的训练集,但是这种差异和装袋一样是靠运气,所不同的是它是前一个基学习器的误差的函数。提升针对一个特定问题的实际性能显然依赖于数据和基学习器。为此,需要有充足的训练数据,并且基学习器应当是弱的但又不是太弱,而且提升对噪声和离群点尤其敏感。

AdaBoost 已经被推广到回归:由 Avnimelech 和 Intrator(1997)提出的一种直截了当的方法是,检查预测误差是否大于某个阈值,如果是则将其标记为错误,然后使用正规的 AdaBoost。在另一个版本中(Drucker1997),抽取概率根据误差量进行修改,使得前一个基学习器预测误差较大的实例,在下一个基学习器的训练中有较大的概率被抽取。最后用加权平均或中值来组合这些基学习器的预测结果。

433

17.8 重温混合专家模型

在投票中,权重 w_j 在输入空间上是固定的。在混合专家模型构架中(参见 12.8 节),存在一个门网络,其输出取作投票的权重。因而这一构架可以被看作是一种投票方法,其中投票依赖于输入,而且可能因输入不同而有所不同。混合专家模型使用的竞争学习算法局部化了基学习器,使得每个基学习器变成输入空间的一个不同领域的专家,并且其权重 $w_j(x)$ 在其专长的领域中接近于 1。最终的输出与投票一样是加权平均

$$y = \sum_{j=1}^L w_j(x) d_j \quad (17.7)$$

不同之处在于基学习器和权重二者均是输入的一个函数(见图 17-3)。

434

Jacobs(1997)显示在混合专家模型构架中,专家是有偏的,但是负相关的。随着训练的进行,偏倚降低而专家的方差增加,但与此同时,随着专家局部化于输入空间的不同部分,它们的协方差为负并且越来越小。根据式(17.6),这将降低总体方差,进而降低误差。在 12.8 节,我们讨论了专家和门网络均为线性函数的情况,但是非线性方法同样可以用于二者。这将降低专家的偏倚,但是有增加专家方差和过分拟合的风险。

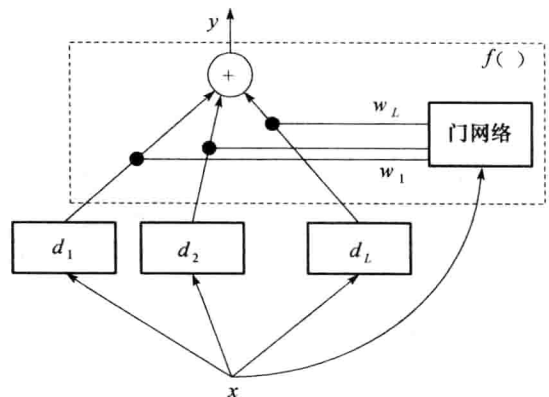


图 17-3 混合专家模型是一种投票方法,其中,像由门网络给出的那样,投票是输入的函数。组合系统 f 也包含这种迭通系统

在动态分类器选择(dynamical classifier selection)中,类似于混合专家模型的门网络,首先由系统取检验输入并评估基学习器在输入周围的竞争力。然后,它挑选最有竞争力的基学习器来产生输出,并将该输出作为总体输出。Woods、Kegelmeyer 和 Bowyer(1977)找出了检验输入的 k 个最近的训练点,考察基分类器在它们上的准确率,从中选择性能最好的一个,只需要用选定的基分类器对该检验输入求值。为了

降低方差, 以更多计算为代价, 我们可以取几个竞争的基分类器的投票, 而不是只用一个。

注意, 在这种模式下, 需要确保对于输入空间的任何区域都存在一个竞争的基分类器。这意味着在基分类器之间存在输入空间的学习划分。这是混合专家模型一个很好的性质, 即做选择的选通模型和它选择的专家基学习器以耦合方式训练。这种动态学习器选择算法有一个显而易见的回归版本(习题 5)。

17.9 层叠泛化

层叠泛化(stacked generalization)是 Wolpert(1992)提出的一种扩展的投票方法, 其中基学习器的输出组合方式不必是线性的, 而是通过一个组合器系统 $f(\cdot | \Phi)$ 。组合器是另一个学习器, 其参数 Φ 也要训练(参见图 17-4):

$$y = f(d_1, d_2, \dots, d_L | \Phi) \quad (17.8)$$

当基学习器给出某种输出组合时, 组合器学习什么是正确的输出。我们不能在训练数据上训练组合器函数, 因为基学习器可能在记忆训练数据; 组合器系统应当真正的学习基学习器是如何犯错(产生误差)的。因此组合器应当在基学习器训练时没有出现的数据上训练。

如果 $f(\cdot | w_1, \dots, w_L)$ 是线性模型, 其约束为 $w_i \geq 0$, $\sum_j w_j = 1$, 则最佳权重可通过受约束的回归来获得。但是请注意, 对组合器函数没有限制, 并且不像投票, 组合可以是非线性的。例如, $f(\cdot)$ 可以是一个多层感知器, Φ 是其连接权重。基学习器 d_j 的输出定义了一个新的 L 维空间, 在该空间组合器函数学习输出的判别式/回归函数。

在层叠泛化中, 我们希望基学习器尽可能不同, 使得它们可以相互补充。为此, 最好每个基学习器都基于不同的学习算法。如果我们组合可以产生连续输出(如后验概率)分类器, 则最好是组合这些连续输出, 而不是组合硬决策。

当我们将一个训练过的组合器(如层叠中的组合器)与诸如投票中的固定规则进行比较时, 我们看到二者都有各自的优点: 训练过的规则更灵活并且可能具有更小的偏倚, 但是它增加了额外的参数, 有引入方差的风险, 并且需要更多的时间和数据进行训练。还要注意, 在进行层叠之前, 不必对分类器的输出规范化。

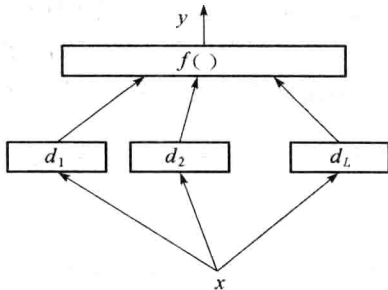


图 17-4 在层叠泛化中, 组合器是另一个学习器, 并且不必像投票一样是线性的

17.10 调整系统

模型组合并非总是能保证降低误差的神奇方法; 基学习器应该是有差异的和准确的, 即它们应该提供有用的信息。如果一个基学习器不能提高准确率, 则可以丢弃它。此外, 如果两个基学习器是高度相关的, 则其中一个不需要。注意, 一个不准确的学习器也可能使准确率恶化; 例如, 多数表决假定超过一半的分类器对一个输入是准确的。因此, 给定一组候选基学习器, 那么使用所有的基学习器可能不是一个好主意, 而选择一个子集可能更好。这意味着选择子集不仅降低了复杂性, 而且也能提高准确率。

从基学习器的系综中选择一个子集类似于输入特征选择, 并且系综选择(ensemble selec-

tion)的方法可能也与特征选择相同。我们可以用向前/增量/增长方法,即每次迭代,从候选基学习器的集合中选择一个最能提高准确率的基学习器,添加到系综中。我们可以用向后/减量/剪枝方法,在每次迭代中从系综中删除一个因为其缺失而能够最大提高性能的基学习器。或者,我们可以使用浮动方法,添加和移去。组合策略可以是固定的规则,如投票,也可以是训练的层叠。这种选择将排除不准确的学习器,排除那些差异性不够或相关的学习器(Caruana等2004, Ruta和Gabrys 2005)。不同的学习器可以使用不同的表示,而这种方法也可以选择最佳的互补表示(Demir和Alpaydin 2005)。

实际上,正如在层叠中一样,如果我们将组合看作一个取基学习器输出为输入的学习器,则我们此处的目标是输入维度归约,而维度归约在第6章讨论过。一种可能的方法仍然是特征选择,丢弃那些不提供信息的输入,而保留那些有用的输入。在系综方法中,正如我们前面讨论的,这对应于从基学习器的系综中选择一个子集。注意,如果我们使用决策树作为组合器,则它既充当选择器,又充当组合器(Ulas等2009)。

第二种可能的方法是特征提取,其目的是从基学习器的输出空间到较低维的空间,在低维空间我们删除不必要的输入,并且也消除相关性。Merz(1999)提出了SCANN算法,即在基分类器的输出上使用对应分析(主成分分析(参见6.3节)的一种变形),并使用最近均值分类器组合它们。实际上,我们在第6章讨论的任意线性或非线性特征提取方法都可以使用,并且它的(最好是连续的)输出可以提供给任意学习器。因此,利用每个 K 个输出的 L 个学习器,我们由 $K \cdot L$ 维空间映射到较低维的新空间,这些“特征学习器”的不相关空间,在这个空间我们训练组合器(使用训练基学习器和维归约时未使用的独立数据集)。

不是彻底地丢弃或保留系综的一个子集,这种方法使用所有的基学习器,因而使用了所有的信息,但并不降低复杂性。

17.11 级联

级联分类器的思想是使用一个基学习器 d_j 的序列,按照空间和时间复杂度或它们使用的数据表示的代价对其进行排序,使得 d_{j+1} 的代价高于 d_j (Kaynak和Alpaydin2000)。级联(cascading)是一种多级方法,并且只有在所有前驱学习器 $d_k(k < j)$ 都不足够确信时才使用 d_j (见图17-5)。为此,与每个学习器相关联的是一个置信度 w_j ,当有 $w_j > \theta_j$ 时我们说 d_j 对其输出是确信的并且其结果可用,其中 $1/K < \theta_j \leq \theta_{j+1} < 1$ 是置信度阈值。在分类中,置信度函数被设置为最高的后验: $w_j \equiv \max_i d_{ji}$;这正是用于拒绝的策略(3.3节)。

如果所有前驱学习器的结果均不够确信,我们才使用学习器 d_j :

$$y_i = d_{ji} \quad \text{如果 } w_j > \theta_j \text{ 并且对于所有 } k < j \text{ 有 } w_k < \theta_k \quad (17.9)$$

给定一个训练集,我们从 $j=1$ 开始训练 d_j 。然后我们从另外一个验证集中找出所有使 d_j 不够好(不确信)的实例,将其组成 d_{j+1} 的训练集。注意,和AdaBoost不同的是,这里我们不仅选取在前一个基学习器上误分类的实例,同时还选取其不自信的实例。这包括误分类的实例以及后验概率不够高的实例;这些实例位于边界正确一侧,但是它们与判别式之间的距离(即边缘)不够大。

级联的思想是:在初期使用简单的分类器处理大多数实例,而更为复杂的分类器仅用于少数实例,因此并不显著增加总体复杂度。这正和类似于投票的多专家方法相反,这些方法

436

437

438

中,所有基学习器为每个实例产生输出。如果问题空间比较复杂,几个每一级的复杂性递增的基学习器可能级联。为了不增加基分类器的个数,少数没有被任何基分类器覆盖的实例将被原样保留,并通过一个非参数分类器(如 k -NN)来处理。

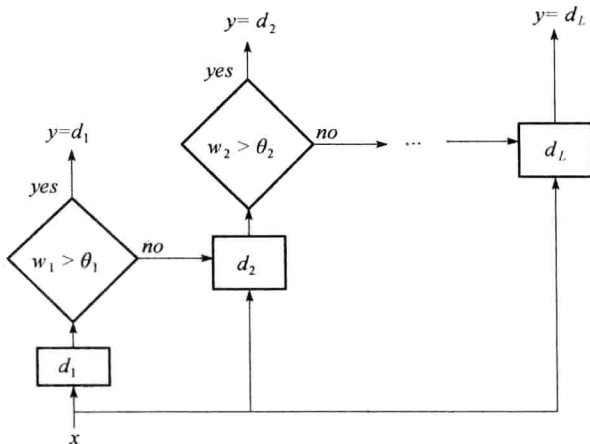


图 17-5 级联是一个多级方法,其中使用一个分类器序列,并且仅当前驱分类器不够好时才使用下一个分类器

级联的归纳偏倚是类可以通过复杂度递增的少量“规则”来解释,并存在一个没有被这些规则覆盖的小的“异常”集合。这些规则通过简单的基学习器来实现,例如,复杂度递增的感知器,学习在整个输入空间上有效的一般规则。异常是局部实例,最好由非参数模型处理。

因此,级联位于参数和非参数分类两个极端之间。前者(例如线性模型)寻找覆盖所有实例的单个规则。而非参数模型(如 k -NN)存储所有的实例而不产生任何解释它们的简单规则。级联产生一个(或多个)规则,以尽可能低代价地解释大部分实例,并将其余实例作为异常存储。这在很多学习应用中是有道理的。例如,在多数情况下,英语动词的过去式是在其后加“d”或“ed”;但也存在不规则动词并不符合这一规则的情况,例如“go/went”。

17.12 注释

组合学习器的思想是将复杂的任务划分为较简单的子任务,这些子任务可以由分别训练的基学习器处理。每个基学习器有其自己的子任务。如果我们用一个大的学习器包含所有的基学习器,则会有过分拟合的风险。例如,考虑取三个多层感知器上的投票,每个感知器具有一个隐藏层。如果我们用线性模型组合其输出,将它们组合在一起,则我们有一个大的、具有两个隐藏层的多层感知器。如果我们用全部样本来训练这个大模型,则很可能产生过分拟合。而当我们分别对三个多层感知器训练时,比如用 ECOC、装袋等,就如同为这个大的多层感知器的第二层隐藏节点定义了所需要的输出。这就为整体的学习器应对什么进行学习附加了约束,进而简化了学习任务。

组合的一个缺点是组合系统不是可解释的。例如,即使决策树是可解释的,装袋的或提升的决策树也不是可解释的。具有权重如 $-1/0/+1$ 的纠错码允许某种形式的可解释性。Mayoraz 和 Moreira(1997)讨论了用于学习纠错输出码的递增方法,其中基学习器在需要时添加。Allwein、Schapire 和 Singer(2000)讨论了将多类目标问题用两类目标问题编码的各种方法。Alpaydm 和 Mayoraz(1999)考虑了 ECOC 的应用,其中对线性基学习器组合得到非线性

判别式，他们还提出了从数据中学习 ECOC 矩阵的方法。

最早也是最直观的方法就是投票。Kittler 等(1998)回顾了固定规则，还讨论了一种组合多种数据表示的应用。该任务是使用三种表示进行人脸识别：正面人脸图像、人脸轮廓图像和声音。投票模型的误差率低于使用单一表示的误差率。Alimoglu 和 Alapaydin(1997)给出了另一种应用，其中为了改善手写数字的识别，对两种信息源进行了组合：一种是数字在触摸书写板书写时的笔移动时态数据，另一种是数字书写后的静态二维位图图像。在这个应用中，使用其中一种数据表示的两个分类器误差率约为 5%，但是通过组合误差率降至 3%。应用研究还表明关键在于设计互补的学习器和/或数据表示，学习器的组合方式并非是至关重要的。

组合不同模态广泛用于生物测定学，这个学科的目标是使用指纹、签名、面部等不同的输入源进行身份认证。在这种情况下，不同的分类器分别使用不同的模态，而它们的预测被组合，这提高了准确率并使得欺骗更困难。

当我们有来自不同表示或模态的不同数据源时，Noble(2004)区分了三种类型的组合策略：

- 在早期集成(early integration)中，所有这些输入都串接起来，形成一个向量，然后提供给一个分类器。前面我们讨论过为什么这不是一个很好的主意。
- 后期集成(late integration)是本章提倡的；在后期集成中，不同的输入提供给单独的分类器，然后通过投票、层叠或我们讨论过的任何方法组合它们的输出。
- 我们在第 13 章讨论的核算法容许一种不同的集成方法，Noble(2004)称其为中期集成(intermediate integration)，因为它介于早期集成与后期集成之间。这是一种多核学习(multiple kernel learning)方法(参见 13.8 节)，其中只有一个核机器分类器，它对不同的输入使用多个核函数，并且组合不像前期集成那样在输入空间进行，也不像后期集成那样在决策空间中进行，而是在定义核函数的基函数空间中进行。对于不同的源，存在由其核函数计算的不同的相似性概念，而分类器累积并使用它们。

Jacobs(1995)表明 L 个依赖的专家和 L' 个独立的专家同样有价值，其中 $L' \leq L$ 。在某些环境下，投票模型与贝叶斯技术产生相同的结果(Jacobs1995)。式(17.3)的先验概率因而可以由超参数的分布来建模，并在理想情况下，应在整个模型参数空间进行集成。这种方法在实际中并非总是可行的，并且需要求助于近似或抽样。随着贝叶斯统计的进展，这些超贝叶斯技术可能会在不远的将来变得越来越重要。

自 20 世纪 90 年代初以来，组合多学习器就已经成为机器学习领域中的一个流行课题，从那时起研究一直在进行。Kuncheva(2004)讨论了分类器组合的方方面面；该书还包含讨论多个聚类结果的组合的章节。

AdaBoost 提升的决策树被认为是最好的机器学习算法之一。同时还存在 AdaBoost 的其他版本，其中下一个基学习器在前一个学习器的残余之上进行训练(Hastie, Tibshirani 和 Friedman, 2001)。最近，人们注意到系综并非总能提高准确率，并且研究者们开始关注好的系综的标准，以及如何得到好的系综。关于系综中差异性的作用的综述在 Kuncheva 2005 中。

17.13 习题

1. 如果每个基学习器是独立同分布的并且正确的概率 $p > 1/2$ ，那么 L 个分类器上的一个多数

441
442

表决给出正确答案的概率是什么?

2. 在装袋中, 为了产生 L 个训练集, 以 L 折交叉确认来替代自助法的效果如何?
3. 提出一个学习纠错输出码的增量算法, 其中新的二类问题在需要时添加, 以便更好地解决多类问题。
4. 在混合专家模型中, 我们可以让不同的专家使用不同的输入表示。在这种情况下, 如何设计门网络?
5. 提出一种动态回归器选择算法。
6. 使用线性感知器作为组合器函数的话, 投票和层叠的区别是什么?
7. 在级联中, 为什么要求 $\theta_{j+1} \geq \theta_j$?
8. 为了能够对回归使用级联, 在检验时, 回归器应该能够表示对其输出是否有把握。如何实现这一点?
9. 如何组合多个聚类解的结果?
10. 在 17.10 节, 我们讨论过如果使用决策树作为层叠中的组合器, 则决策树充当选择器和组合器。还有什么其他优点和缺点?

17.14 参考文献

- Alimoğlu, F., and E. Alpaydın. 1997. "Combining Multiple Representations and Classifiers for Pen-Based Handwritten Digit Recognition." In *Fourth International Conference on Document Analysis and Recognition*, 637-640. Los Alamitos, CA: IEEE Computer Society.
- Allwein, E. L., R. E. Schapire, and Y. Singer. 2000. "Reducing Multiclass to Binary: A Unifying Approach for Margin Classifiers." *Journal of Machine Learning Research* 1: 113-141.
- Alpaydın, E. 1997. "Voting over Multiple Condensed Nearest Neighbors." *Artificial Intelligence Review* 11: 115-132.
- Alpaydın, E., and E. Mayoraz. 1999. "Learning Error-Correcting Output Codes from Data." In *Ninth International Conference on Artificial Neural Networks*, 743-748. London: IEE Press.
- Avnimelech, R., and N. Intrator. 1997. "Boosting Regression Estimators." *Neural Computation* 11: 499-520.
- Breiman, L. 1996. "Bagging Predictors." *Machine Learning* 26: 123-140.
- Caruana, R., A. Niculescu-Mizil, G. Crew, and A. Ksikes. 2004. "Ensemble Selection from Libraries of Models." In *Twenty-First International Conference on Machine Learning*, ed. C. E. Brodley, 137-144. New York: ACM.
- Demir, C., and E. Alpaydın. 2005. "Cost-Conscious Classifier Ensembles." *Pattern Recognition Letters* 26: 2206-2214.
- Dietterich, T. G., and G. Bakiri. 1995. "Solving Multiclass Learning Problems via Error-Correcting Output Codes." *Journal of Artificial Intelligence Research* 2: 263-286.
- Drucker, H. 1997. "Improving Regressors using Boosting Techniques." In *Fourteenth International Conference on Machine Learning*, ed. D. H. Fisher, 107-115. San Mateo, CA: Morgan Kaufmann.

- Drucker, H., C. Cortes, L. D. Jackel, Y. Le Cun, and V. Vapnik. 1994. "Boosting and Other Ensemble Methods." *Neural Computation* 6: 1289-1301.
- Freund, Y., and R. E. Schapire. 1996. "Experiments with a New Boosting Algorithm." In *Thirteenth International Conference on Machine Learning*, ed. L. Saitta, 148-156. San Mateo, CA: Morgan Kaufmann.
- Hansen, L. K., and P. Salamon. 1990. "Neural Network Ensembles." *IEEE Transactions on Pattern Analysis and Machine Intelligence* 12: 993-1001.
- Hastie, T., R. Tibshirani, and J. Friedman. 2001. *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. New York: Springer.
- Ho, T. K. 1998. "The Random Subspace Method for Constructing Decision Forests." *IEEE Transactions on Pattern Analysis and Machine Intelligence* 20: 832-844.
- Jacobs, R. A. 1995. "Methods for Combining Experts' Probability Assessments." *Neural Computation* 7: 867-888.
- Jacobs, R. A. 1997. "Bias/Variance Analyses for Mixtures-of-Experts Architectures." *Neural Computation* 9: 369-383.
- Jain, A., K. Nandakumar, and A. Ross. 2005. "Score Normalization in Multimodal Biometric Systems." *Pattern Recognition* 38: 2270-2285.
- Kaynak, C., and E. Alpaydın. 2000. "MultiStage Cascading of Multiple Classifiers: One Man's Noise is Another Man's Data." In *Seventeenth International Conference on Machine Learning*, ed. P. Langley, 455-462. San Francisco: Morgan Kaufmann.
- Kittler, J., M. Hatef, R. P. W. Duin, and J. Matas. 1998. "On Combining Classifiers." *IEEE Transactions on Pattern Analysis and Machine Intelligence* 20: 226-239.
- Kuncheva, L. I. 2004. *Combining Pattern Classifiers: Methods and Algorithms*. Hoboken, NJ: Wiley.
- Kuncheva, L. I. 2005. Special issue on Diversity in Multiple Classifier Systems. *Information Fusion* 6: 1-115.
- Mayoraz, E., and M. Moreira. 1997. "On the Decomposition of Polychotomies into Dichotomies." In *Fourteenth International Conference on Machine Learning*, ed. D. H. Fisher, 219-226. San Mateo, CA: Morgan Kaufmann.
- Merz, C. J. 1999. "Using Correspondence Analysis to Combine Classifiers." *Machine Learning* 36: 33-58.
- Noble, W. S. 2004. "Support Vector Machine Applications in Computational Biology." In *Kernel Methods in Computational Biology*, ed. B. Schölkopf, K. Tsuda, and J.-P. Vert, 71-92. Cambridge, MA: MIT Press.
- Perrone, M. P. 1993. "Improving Regression Estimation: Averaging Methods for Variance Reduction with Extensions to General Convex Measure." Ph.D. thesis, Brown University.
- Ruta, D., and B. Gabrys. 2005. "Classifier Selection for Majority Voting." *Information Fusion* 6: 63-81.
- Schapire, R. E. 1990. "The Strength of Weak Learnability." *Machine Learning* 5: 197-227.

- Schapire, R. E., Y. Freund, P. Bartlett, and W. S. Lee. 1998. "Boosting the Margin: A New Explanation for the Effectiveness of Voting Methods." *Annals of Statistics* 26: 1651-1686.
- Ulaş, A., M. Semerci, O. T. Yıldız, and E. Alpaydın. 2009. "Incremental Construction of Classifier and Discriminant Ensembles." *Information Sciences* 179: 1298-1318.
- Wolpert, D. H. 1992. "Stacked Generalization." *Neural Networks* 5: 241-259.
- Woods, K., W. P. Kegelmeyer Jr., and K. Bowyer. 1997. "Combination of Multiple Classifiers Using Local Accuracy Estimates." *IEEE Transactions on Pattern Analysis and Machine Intelligence* 19: 405-410.

增强学习

在增强学习中，学习器是一个制定决策的智能主体。智能主体在其所处的环境中执行一些动作并根据其试图解决一个问题所执行的动作而获得奖励(或惩罚)。经过反复尝试运行，学习程序应当可以学习得到最优策略，即一个最大化总体奖励的动作序列。

18.1 引言

假设我们要构建一个学习下国际象棋的机器。在这种情况下，我们不能使用监督学习，原因有二：首先，请一位国际象棋老师带领我们遍历许多棋局并告诉我们每个位置的最佳棋步的代价非常昂贵；其次，在很多情况下，根本就没有最佳棋步，一个棋步的好坏依赖于其后的多个棋步。单一的棋步并不算数；而如果经过一个棋步序列我们赢得了比赛，则该棋步序列才是好的。而整个过程唯一的反馈是在最后我们赢得或是输掉棋局时才产生。

另一个例子是置于迷宫中的机器人。机器人按照四个罗盘方向之一进行移动，并进行一系列的移动到达迷宫出口。只要机器人在迷宫中，就不存在反馈，并且机器人尝试各种移动，直至到达出口，只有这时它才得到一个奖励。在这种情况下，机器人不存在对手，但是我们可能更偏好更短的(到达出口)的路径，这意味着我们是在和时间比赛。

这两个应用有多个共同点：存在一个称为智能主体(agent)的决策者，并置其于某一环境(environment)中(见图 18-1)。在国际象棋的例子中，棋手是决策者而环境是棋盘；在第二个例子中，迷宫是机器人的环境。在任何时候，环境总是处于某种状态(state)，该状态来自于一组可能的状态之一，例如，棋盘的布局状态，机器人在迷宫中的位置。决策者可以做一组可能的动作：棋盘上棋子的合法移动，机器人沿着可能的方向移动而不会撞墙等。一旦选择并做了某一动作，状态就随之改变。问题的解决需要执行一系列的动作，之后我们才得到反馈，反馈以极少发生的奖励(reward)的形式给出，通常只有在完整的动作序列执行完毕才发生。奖励对问题进行了定义，同时是构建一个会学习的(learning)智能主体所必需的。会学习的智能主体学习解决一个问题的最佳步骤，其中“最佳”是以获取最大累积奖励的动作序列来量化。以上就是增强学习(reinforcement learning)的背景。

447

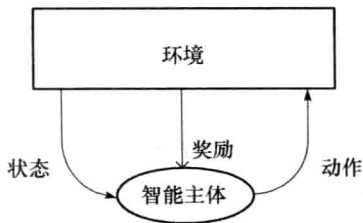


图 18-1 智能主体和环境进行交互。在环境的任意一个状态，智能主体执行一个改变环境状态的动作并获得一个奖励

增强学习与之前讨论的各种学习方法的不同之处在于以下几个方面：它称之为“和批评者一起学习”，而与之前和老师一起学习的监督学习方法相反。批评者(critic)不同于老师之处在于他并不告诉我们做什么，而仅仅告诉我们之前所做的怎么样；批评者永远不会提前提供信息。批评者提供的反馈极少，并且当他提供时，也是事后提供。这就导致了信度分配(credit assignment)问题：在执行若干动作并获得奖励后，我们希望对之前所执行的单个动作进行评估并找到可以引领我们赢得奖励的那些动作，以便对其记录并在之后使用。正如我们即将看到的，一个增强学习程序所做的是为中间状态或动作产生一个内部值(internal value)，来表明这些状态或动作在引领我们达到目标并获取真正的奖励方面有多好。一旦学习到这样的内部奖励机制，智能主体就可以只执行最大化内部奖励的局部的动作。

448

问题的解决需要执行一个动作序列，而从这一角度，我们可想到第 15 章讨论的马尔可夫模型。事实上，我们使用马尔可夫决策过程来对智能主体建模。不同之处在于，对于马尔可夫模型，存在一个外部过程来产生一个我们对其观测和建模的信号序列，如语音。而在增强学习中，产生动作序列的是智能主体。之前，我们还区别可观测的和隐藏的马尔可夫模型，分别对应系统状态是可观测的或是隐藏的(并且也需要推断)。类似地，有时候我们使用一个部分可观测的马尔可夫决策过程来建模，其中智能主体不是确切地知道其所处的状态，而是需要通过使用传感器的观测以某种非确定性对其进行推断。例如，机器人在房间中移动时，机器人可能不知道其在房间中的确切位置，也不知道障碍物和目标的确切位置，而是通过一个照相机提供的有限图像来做决策。

18.2 单状态情况： K 臂赌博机问题

我们从一个简单的例子开始。 K 臂赌博机是一种假想的具有 K 个手柄的老虎机。可做的动作是选择并拉下其中的一个手柄，而由此所赢取的一定量的钱就是和这个手柄(动作)相关联的奖励。任务是决定拉下哪个手柄，以便得到最大奖励。这是一个分类问题，其中我们选择 K 个手柄中的一个。如果是监督学习，则老师会告诉我们正确的类，即产生最大收益的类。而在增强学习中，我们只能尝试不同的手柄并记录其中最好的。这是一个简化的增强学习问题，因为只有一个状态，或者说只有一个老虎机，而我们只需要确定所执行的动作。另一个称其为简化问题的原因是我们在一个动作之后立即得到一个奖励；奖励并没有被延迟，因此在动作之后可以立即看到其价值。

假设 $Q(a)$ 是动作 a 的价值。最初，对所有 a 都有 $Q(a) = 0$ 。当我们尝试执行动作 a 时，我们获得一个奖励 $r_a \geq 0$ 。如果奖励是确定性的，拉下手柄 a 我们总是获得相同的奖励 r_a ，并且在这种情况下，我们可以简单地令 $Q(a) = r_a$ 。如果我们想充分利用已有发现，一旦我们发现一个动作 a 具有 $Q(a) > 0$ ，我们就可以不停地选择它并在每次拉下手柄时获得 r_a 的奖励。然而，很可能还存在另一个具有更高奖励的手柄，因此我们还需要进行探索。

449

我们可以选择不同的动作并对所有的 a 存储 $Q(a)$ 。只要我们想充分利用已有发现，我们可以选择具有最高价值的动作，即

$$\text{选择 } a^*, \quad \text{如果 } Q(a^*) = \max_a Q(a) \quad (18.1)$$

如果奖励并非确定的而是随机的，则在选择相同的动作时我们每次获得不同的奖励。奖励量被概率分布 $p(r|a)$ 定义。在这种情况下，我们定义 $Q_t(a)$ 作为在时刻 t 时执行动作 a 的

价值的估计。它是在时刻 t 之前所有执行动作 a 所获奖励的平均值。一种在线更新方法可定义如下：

$$Q_{t+1}(a) \leftarrow Q_t(a) + \eta[r_{t+1}(a) - Q_t(a)] \quad (18.2)$$

其中 $r_{t+1}(a)$ 是在 $(t+1)$ 时刻执行动作 a 之后所获得的奖励。

注意式(18.2)正是我们在之前章节中多次使用的 delta 规则(delta rule)： η 是学习因子(为了收敛而随时间逐渐降低)， r_{t+1} 是期望输出，而 $Q_t(a)$ 是当前的预测。 $Q_{t+1}(a)$ 是在 $t+1$ 时刻的动作 a 的期望值，并且随着 t 的增加收敛到 $p(r|a)$ 的均值。

完整的增强学习问题从以下几个方面推广了这种简单情况：首先，我们有多状态。这相当于同时存在具有多个不同奖励概率 $p(r|s_i, a_j)$ 的老虎机，而我们需要对 $Q(s_i, a_j)$ ，即在状态 s_i 时执行动作 a_j 的价值进行学习；其次，动作不仅影响获得的奖励而且影响下一状态，并且我们从一个状态转移到另一个状态；第三，奖励被延迟，而我们需要能够从延迟的奖励值估计立即值。

18.3 增强学习基础

学习的决策者称为智能主体(agent)。智能主体与环境(environment)之间进行交互。环境包含了除智能主体之外的所有东西。智能主体具有感知器，用来决定其所处的状态(state)。当智能主体执行一个动作时，环境提供一个奖励(reward)。时间被离散化为 $t=0, 1, 2, \dots$ ，并且 $s_t \in S$ 表示智能主体在时刻 t 的状态，其中 S 是所有可能的状态集合。 $a_t \in \mathcal{A}(s_t)$ 表示智能主体在时刻 t 所执行的动作，其中 $\mathcal{A}(s_t)$ 是在状态 s_t 时所有可能执行的动作集合。当处于状态 s_t 的智能主体执行动作 a_t 时，时钟嘀嗒，接收到奖励 $r_{t+1} \in \mathbb{R}$ ，并且智能主体转移到下一个状态 s_{t+1} 。整个问题通过马尔可夫决策过程(Markov decision process, MDP)来建模。奖励和下一状态分别采样于它们相应的概率分布 $p(r_{t+1}|s_t, a_t)$ 和 $P(s_{t+1}|s_t, a_t)$ 。注意，我们所具有的是一个马尔可夫(Markov)系统，其中下一时刻的状态和奖励仅仅依赖于当前状态和动作。在一些应用中，奖励和下一状态是确定的，并且对某个状态和所执行的动作，存在一个可能的奖励值和下一状态。

[450]

依赖于应用，某一状态可能被指定为初始状态，而在一些应用中，也存在一个停止搜索的吸收终止(目标)状态；所有在终止状态执行的动作都以概率 1 将状态转移到自身并且没有任何奖励。从初始状态到终止状态的动作序列称为一个片段(episode)，或一次试验(trial)。

策略(policy) π 定义了智能主体的行为并且是从环境状态到动作之间的一个映射： $\pi: S \rightarrow \mathcal{A}$ 。策略定义了在任何状态 s_t 可以执行的动作： $a_t = \pi(s_t)$ 。策略 π 的价值 $V^\pi(s_t)$ 是从状态 s_t 开始，遵循该策略的智能主体所获得的期望累积奖励。

在有限视野(finite-horizon)或片段(episodic)模型中，智能主体试图最大化下 T 个步骤的期望奖励：

$$V^\pi(s_t) = E[r_{t+1} + r_{t+2} + \dots + r_{t+T}] = E\left[\sum_{i=1}^T r_{t+i}\right] \quad (18.3)$$

某些任务是连续的，并且不存在预先固定的关于片段的限制。在无限视野(infinite-horizon)模型中，不存在序列长度的限制，但是未来的奖励将被打折扣：

$$V^\pi(s_t) = E[r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots] = E\left[\sum_{i=1}^{\infty} \gamma^{i-1} r_{t+i}\right] \quad (18.4)$$

其中 $0 \leq \gamma < 1$ 是折扣率 (discount rate), 保证所返回的奖励是有限的。如果 $\gamma = 0$, 则只有立即的奖励算数。随着 γ 趋向于 1, 处于未来的奖励将更多地被计算在内, 而这时我们说智能主体变得更有远见了。 γ 是小于 1 的, 因为对于解决问题的动作序列总是会有一个时间上的限制, 而且智能主体也许是一个靠电池运作的机器人。因此我们喜欢更早获得奖励而非更晚, 因为我们不确定智能主体可以运行多久。

对每个策略 π , 存在其价值 $V^\pi(s_t)$, 我们想要找到最优策略 (optimal policy) π^* 使得

$$V^*(s_t) = \max_{\pi} V^\pi(s_t), \forall s_t \quad (18.5)$$

在某些应用中, 例如在控制中, 我们更希望处理成对的状态 - 动作值 $Q(s_t, a_t)$, 而不是简单的状态值 $V(s_t)$ 。 $V(s_t)$ 表示智能主体处于状态 s_t 的价值, 而 $Q(s_t, a_t)$ 表示当处于状态 s_t 时执行动作 a_t 的价值。我们定义 $Q^*(s_t, a_t)$ 为处于 s_t 状态时执行动作 a_t 并在其后遵循最优策略的期望累积奖励。状态的价值等于其上可采取的最优动作的价值:

$$\begin{aligned} V^*(s_t) &= \max_{a_t} Q^*(s_t, a_t) \\ &= \max_{a_t} E \left[\sum_{i=1}^{\infty} \gamma^{i-1} r_{t+i} \right] \\ &= \max_{a_t} E \left[r_{t+1} + \gamma \sum_{i=1}^{\infty} \gamma^{i-1} r_{t+i+1} \right] \\ &= \max_{a_t} E [r_{t+1} + \gamma V^*(s_{t+1})] \\ V^*(s_t) &= \max_{a_t} \left(E[r_{t+1}] + \gamma \sum_{s_{t+1}} p(s_{t+1} | s_t, a_t) V^*(s_{t+1}) \right) \end{aligned} \quad (18.6)$$

对于每一个可能的下一状态 s_{t+1} , 我们以概率 $P(s_{t+1} | s_t, a_t)$ 转移到 s_{t+1} 并自此遵循最优策略, 所得的期望累积奖励是 $V^*(s_{t+1})$ 。我们在所有可能的下一状态上求和, 并且打折扣, 因为它晚一个时间步。加上立即期望奖励, 我们得到动作 a_t 的总体期望累积奖励。最后我们选择所有动作中最好的一个。式 (18.6) 称为 Bellman 公式 (Bellman's equation) (Bellman 1957)。类似地, 我们还可以有

$$Q^*(s_t, a_t) = E[r_{t+1}] + \gamma \sum_{s_{t+1}} p(s_{t+1} | s_t, a_t) \max_{a_{t+1}} Q^*(s_{t+1}, a_{t+1}) \quad (18.7)$$

一旦获得了 $Q^*(s_t, a_t)$ 的值, 我们就可以定义策略 π 为执行动作 a_t^* , 它在所有 $Q^*(s_t, a_t)$ 中具有最大值:

$$\pi^*(s_t): \text{选择 } a_t^*, \text{ 其中 } Q^*(s_t, a_t^*) = \max_{a_t} Q^*(s_t, a_t) \quad (18.8)$$

这意味着只要我们获得所有 $Q^*(s_t, a_t)$ 的值, 那么在每个局部步骤中使用贪心搜索, 我们就可以得到一个最优的步骤序列, 该序列最大化累积 (cumulative) 奖励。

18.4 基于模型的学习

我们从基于模型的学习开始, 其中我们完全知晓环境模型的参数 $p(r_{t+1} | s_t, a_t)$ 和 $P(s_{t+1} | s_t, a_t)$ 。在这种情况下, 我们不要进行任何探索就可以使用动态规划直接对最优价值函数和策略求解。最优价值函数是唯一的, 即为式 (18.6) 的解。一旦获得了最优价值函数, 最优策略即为选择最大化下一状态价值的动作:

$$\pi^*(s_t) = \arg \max_{a_t} \left(E[r_{t+1} | s_t, a_t] + \gamma \sum_{s_{t+1} \in S} P(s_{t+1} | s_t, a_t) V^*(s_{t+1}) \right) \quad (18.9)$$

18.4.1 价值迭代

为了找到最优策略，可以使用最优价值函数，并且存在一个称为价值迭代 (value iteration) 的迭代算法，业已证明它收敛于正确的 V^* 值。价值迭代算法的伪代码在图 18-2 中。

```

将  $V(s)$  初始化为任意值
Repeat
  For 所有的  $s \in S$ 
    For 所有的  $a \in \mathcal{A}$ 
       $Q(s, a) \leftarrow E[r | s, a] + \gamma \sum_{s' \in S} P(s' | s, a) V(s')$ 
     $V(s) \leftarrow \max_a Q(s, a)$ 
Until  $V(s)$  收敛
  
```

图 18-2 基于模型学习的价值迭代算法

我们说价值迭代是收敛的，条件为两次迭代之间的最大价值差小于某个阈值 δ ：

$$\max_{s \in S} |V^{(l+1)}(s) - V^{(l)}(s)| < \delta$$

453

其中 l 是迭代计数。由于我们只关心具有最大价值的动作，因此有可能在价值收敛于最优价值之前策略就收敛于最优策略。每次迭代的复杂度是 $\mathcal{O}(|S|^2 |\mathcal{A}|)$ ，但是下一个可能状态数目 $k < |S|$ 很小，因此复杂度降低到 $\mathcal{O}(k |S| |\mathcal{A}|)$ 。

18.4.2 策略迭代

在策略迭代中，我们直接存储和更新策略，而非间接地通过价值迭代寻求最优策略。图 18-3 给出了其伪代码。其思想是从一个策略开始，不断的改进它直到没有改变为止。价值函数可通过求解线性方程来计算。然后检验是否可以通过将这些解考虑在内而改进策略。这一步骤保证了对策略的改进，并且当不再可能继续改进时，可以确保所得策略是最优的。该算法每次迭代的时间复杂度是 $\mathcal{O}(|\mathcal{A}| |S|^2 + |S|^3)$ ，比价值迭代的复杂度高，但是策略迭代比价值迭代需要更少的迭代次数。

```

任意初始化一个策略  $\pi'$ 
Repeat
   $\pi \leftarrow \pi'$ 
  通过解线性方程组，计算使用  $\pi$  的价值
   $V^\pi(s) = E[r | s, \pi(s)] + \gamma \sum_{s' \in S} P(s' | s, \pi(s)) V^\pi(s')$ 
  在每个状态上改进策略
   $\pi'(s) \leftarrow \operatorname{argmax}_a (E[r | s, a] + \gamma \sum_{s' \in S} P(s' | s, a) V^\pi(s'))$ 
Until  $\pi = \pi'$ 
  
```

图 18-3 基于模型学习的策略迭代算法

18.5 时间差分学习

454

模型通过奖励和下一状态概率分布来定义，而且从 18.4 节可以看到，当这些值均已知时，可以使用动态规划来求解最佳策略。然而，这些方法代价很高，并且我们很少具有如此完全的关于环境的知识。增强学习更有趣和更实际的应用是当我们并不知道模型的时候。这时，我们需要对环境进行探索来查询模型。我们首先讨论如何进行探索，而后我们讨论在确定和非确定情况下的无模型学习算法。尽管我们并不假定关于环境模型的全部知识是已知的，但是还是要求模型是固定的。

像我们稍后将要看到的，当我们进行探索并得以看到下一个状态的价值和奖励时，我们利用这一信息来更新当前状态的价值。这些算法称为时间差分 (temporal difference) 算法，因为我们所做的是考察一个状态 (或状态 - 动作对) 的价值的当前估计值与下一状态和所得到奖励的折扣值之间的差。

18.5.1 探索策略

为了对环境进行探索，一种可能性是使用 ε -贪心 (greedy) 搜索，其中我们以概率 ε 在所有可能的动作中均匀、随机地选择一个动作，即进行探索；而以概率 $1 - \varepsilon$ 选择已知的最好动作，即进行利用。我们并不想无限地持续探索，而是一旦进行了足够的探索就开始对其利用；为此，我们以一个较大的 ε 值开始，并逐渐减小它。我们需要确认所采取的策略是软 (soft) 策略，也就是说，在状态 $s \in S$ 执行任意动作 $a \in \mathcal{A}$ 的概率大于 0。

我们可以根据概率进行选择，使用软最大函数将价值转化为概率

$$P(a | s) = \frac{\exp Q(s, a)}{\sum_{b \in \mathcal{A}} \exp Q(s, b)} \quad (18.10)$$

然后根据这些概率对动作进行选择。为了逐渐地从探索向利用进行转移，我们可以使用一个“温度”变量 T ，并定义选择动作 a 的概率为

$$P(a | s) = \frac{\exp[Q(s, a)/T]}{\sum_{b \in \mathcal{A}} \exp[Q(s, b)/T]} \quad (18.11)$$

当 T 很大的时候，所有的概率相等，因而我们进行的是探索。而当 T 很小的时候，更好的动作将受青睐。因此这时的策略是以一个大的 T 值开始并逐渐减小它，这称为退火 (annealing) 过程，在这种情况下就是在时间上从探索平滑地过渡到利用。

455

18.5.2 确定性奖励和动作

在无模型学习中，我们首先讨论较为简单的确定性情况，其中对任意一对状态 - 动作，只有一个奖励和可能的下一状态。在这种情况下，式 (18.7) 简化为

$$Q(s_t, a_t) = r_{t+1} + \gamma \max_{a_{t+1}} Q(s_{t+1}, a_{t+1}) \quad (18.12)$$

而我们简单地将其作为一个赋值来更新 $Q(s_t, a_t)$ 。当在状态 s_t 时, 我们使用之前所见到的各种随机策略之一选择一个动作 a_t , 其返回一个奖励 r_{t+1} 并将状态转移到 s_{t+1} 。然后, 前一个动作的价值更新为

$$\hat{Q}(s_t, a_t) \leftarrow r_{t+1} + \gamma \max_{a_{t+1}} \hat{Q}(s_{t+1}, a_{t+1}) \quad (18.13)$$

其中 Q 之上加帽表示该值为估计值。 $\hat{Q}(s_{t+1}, a_{t+1})$ 是一个稍晚的值, 因此有更高的概率是正确的。我们以 γ 对其进行折扣并加上立即奖励(如果有的话), 并将此作为前一个 $\hat{Q}(s_t, a_t)$ 的新估计。这称为后推(backup), 因为这可以被看作是在下一个时间步骤中对一个动作的价值进行估计, 并“将其后退”用来修改一个当前动作的价值估计。

目前, 我们假定所有的值 $\hat{Q}(s, a)$ 存储于一张表中; 稍后, 我们会讨论当 $|S|$ 和 $|A|$ 很大时如何更为简洁地存储这些信息。

最初所有 $\hat{Q}(s_t, a_t)$ 都为0, 并且作为试验片段的结果及时更新。假设我们有一个状态转移的序列, 并且在每次转移中, 我们都使用式(18.13), 用当前状态-动作对的 Q 值来更新前一对状态-动作的 Q 值的估计。在中间状态, 所有的奖励为0从而价值为0, 因此不进行更新。当到达目标状态的时候, 我们得到奖励 r , 因而可以将前一对状态-动作的 Q 值更新为 γr 。对于这个之前的状态-动作对, 由于其立即奖励为0而来自下一对状态-动作对的贡献又因为晚一步而以 γ 折扣。于是, 在下一个片段中, 如果我们再次到达这个状态, 我们将其前一状态更新为 $\gamma^2 r$, 以此类推。按照这种方式, 经过多个片段之后, 这一信息被后推到更早的状态-动作对。随着我们寻找到具有更高累积奖励的路径, 如更短的路径, Q 值不断递增直到最优值, 而且这些 Q 值绝不会降低(参见图18-4)。

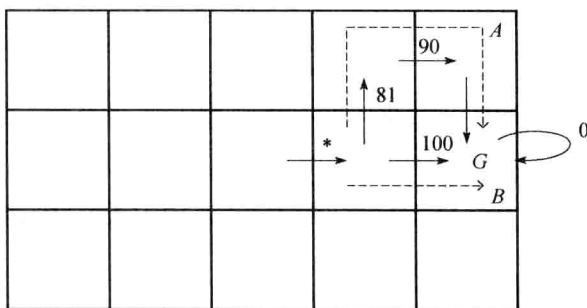


图 18-4 说明 Q 值只增不减的例子。图示是一个确定的网格世界, 其中 G 是目标状态并具有奖励100, 所有其他立即奖励为0并有 $\gamma=0.9$ 。考虑由星号标记的转移的 Q 值, 而且只考虑 A 和 B 两条路径。假设在看到路径 B 之前先看到路径 A , 则有 $\gamma \max(0, 81) = 72.9$ 。如果之后又看到了 B , 则找到了更短的路径, 而 Q 的值变为 $\gamma \max(100, 81) = 90$ 。如果路径 B 在 A 之前被看到, 则 Q 值为 $\gamma \max(100, 0) = 90$ 。于是, 当看到 A 时, Q 的值不变, 因为 $\gamma \max(100, 81) = 90$

注意, 这里我们并不知道奖励或下一状态函数。这些是环境的一部分, 就好像是我们在探索的时候对其进行查询。我们也不对其进行建模, 虽然有此可能性。我们只是原样接受它们并通过估计的价值函数来直接学习最优策略。

18.5.3 非确定性奖励和动作

如果奖励和动作的结果不是确定性的,则我们有一个奖励从中抽样的概率分布 $p(r_{t+1} | s_t, a_t)$, 并且存在一个下一状态所服从的概率分布 $P(s_{t+1} | s_t, a_t)$ 。这些概率分布函数帮助我们对环境不可控制的力量所引发的非确定性进行建模。这些不可控制的力量如国际象棋中的对手, 西洋双陆棋中的骰子, 或者是我们对系统知识的匮乏。例如, 或许我们有一个不完美的机器人, 它有时候会无法按预定的方向前进而产生偏离, 或者比期望的距离走得更近或更远。

在这种情况下, 我们有

$$Q(s_t, a_t) = E[r_{t+1}] + \gamma \sum_{s_{t+1}} P(s_{t+1} | s_t, a_t) \max_{a_{t+1}} Q(s_{t+1}, a_{t+1}) \quad (18.14)$$

在这种情况下, 我们不能进行直接的赋值, 因为对于相同的状态或动作, 我们可能获得不同的奖励或者转移到不同的下一状态。我们所做的是取移动平均。这称之为 Q 学习 (Q learning) 算法:

$$\hat{Q}(s_t, a_t) \leftarrow \hat{Q}(s_t, a_t) + \eta(r_{t+1} + \gamma \max_{a_{t+1}} \hat{Q}(s_{t+1}, a_{t+1}) - \hat{Q}(s_t, a_t)) \quad (18.15)$$

我们将这些 $r_{t+1} + \gamma \max_{a_{t+1}} \hat{Q}(s_{t+1}, a_{t+1})$ 值看作每个 (s_t, a_t) 对的实例的一个样本, 并希望 $\hat{Q}(s_t, a_t)$ 收敛到其均值。与通常一样, 为了收敛, η 的值随时间递减, 并且已经证明该算法收敛于最优的 Q^* 值 (Watkins 和 Dayan 1992)。 Q 学习算法的伪代码见图 18-5。

```

任意初始化所有的  $Q(s, a)$ 
For 所有的片断
    初始化  $s$ 
    Repeat
        使用由  $Q$  导出的策略 (例如  $\epsilon$ -贪心) 选择  $a$ 
        执行动作  $a$ , 观测  $r$  和  $s'$ 
        更新  $Q(s, a)$ :
             $Q(s, a) \leftarrow Q(s, a) + \eta(r + \gamma \max_{a'} Q(s', a') - Q(s, a))$ 
         $s \leftarrow s'$ 
    Until  $s$  是终止状态
    
```

图 18-5 Q 学习, 它是一种离策略时间差分学习算法

我们还可以认为式 (18.15) 的作用是减小当前的 Q 值和一个时间步骤之后的被后推的估计之间的差。这类算法称为时间差分 (temporal difference, TD) 算法 (Sutton 1988)。

这是一种离策略 (off-policy) 方法, 因为该方法使用下一个最优动作的值而不使用策略。在一个在策略 (on-policy) 方法中, 策略还用于确定下一个动作。 Q 学习的在策略版本就是 Sarsa 算法, 其伪代码见图 18-6。我们看到, 在策略的 Sarsa 算法使用从 Q 值推演出的策略来选择下一个动作 a' , 并使用该动作的 Q 值来计算时间差分, 而不是寻找所有可能的下一动作 a' 并选择其中最好的。在策略方法估计一个策略的价值并用它来执行动作。而在离策略方法中, 这些部分都是分离的, 并且用于产生行为的策略称为行为 (behavior) 策略。行为策略事实上可能不同于称为估计 (estimation) 策略的被评估和被改进的策略。

如果采用 GLIE 策略来选择动作, Sarsa 算法以概率 1 收敛到最优策略和状态 - 动作值。GLIE (Greedy in Limit with Infinite Exploration, 使用无限探索的极限贪心) 策略是 1) 所有状态 - 动作对都被无限次访问, 并且 2) 策略收敛到贪心策略的极限 (贪心策略是可设定的, 比如, 使用 ε -贪心策略时设定 $\varepsilon = 1/t$)。

除了 $Q(s, a)$ 之外, 时间差分相同的思想还可以用于学习 $V(s)$ 值。TD 学习 (TD learning) (Sutton1988) 使用如下的更新规则来更新一个状态值:

$$V(s_t) \leftarrow V(s_t) + \eta [r_{t+1} + \gamma V(s_{t+1}) - V(s_t)] \quad (18.16)$$

上式依然是一个 delta 规则, 其中 $r_{t+1} + \gamma V(s_{t+1})$ 是更好的、后一时刻的预测, 而 $V(s_t)$ 是当前的预测。它们之间的差即为时间差分, 而更新是为了减小这个差。更新因子 η 逐渐减小, 因而 TD 确保收敛到最优值函数 $V^*(s)$ 。

```

任意初始化所有的  $Q(s, a)$ 
For 所有的片段
    初始化  $s$ 
    使用由  $Q$  导出的策略 (例如  $\varepsilon$ -贪心) 选择  $a$ 
    Repeat
        执行动作  $a$ , 观测  $r$  和  $s'$ 
        使用由  $Q$  导出的策略 (例如  $\varepsilon$ -贪心) 选择  $a'$ 
        更新  $Q(s, a)$ :
             $Q(s, a) \leftarrow Q(s, a) + \eta (r + \gamma Q(s', a') - Q(s, a))$ 
         $s \leftarrow s', a \leftarrow a'$ 
    Until  $s$  是终止状态
    
```

图 18-6 Sarsa 算法, 它是 Q 学习算法的在策略版本

18.5.4 资格迹

前述算法均为单步算法, 因为时间差分仅用于更新前一个 (状态值或状态 - 动作对的) 值。资格迹 (eligibility traces) 是对以往出现的状态 - 动作对的一个记录, 它使得我们可以实现时间信度分配, 并且还可以更新以往达到的状态 - 动作对的值。我们以 Sarsa 算法学习 Q 值为例来说明这些都是如何完成的。对其进行修改来学习 V 值是直截了当的。

为了存储资格迹, 需要为每个状态 - 动作对关联一个附加的内存变量 $e(s, a)$, 初始化为 0。当状态 - 动作对 (s, a) 被访问时, 也就是说在状态 s 执行了动作 a 时, 其资格被设置为 1; 其他所有状态 - 动作对的资格乘以 $\gamma\lambda$ 。 $0 \leq \lambda \leq 1$ 是迹衰减参数。

$$e_t(s, a) = \begin{cases} 1 & \text{如果 } s = s_t \text{ 并且 } a = a_t \\ \gamma\lambda e_{t-1}(s, a) & \text{否则} \end{cases} \quad (18.17)$$

如果某一状态 - 动作对从未被访问过, 则其资格保持为 0; 如果被访问过, 随着时间流逝和其他状态 - 动作对被访问, 该状态的资格依赖于 γ 和 λ 的值进行衰减 (见图 18-7)。

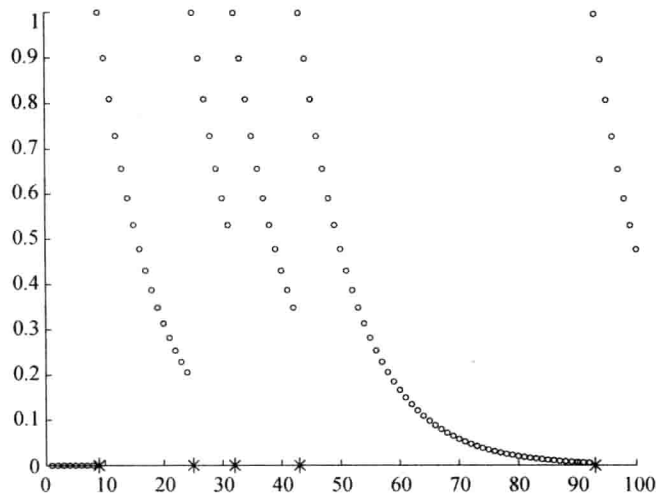


图 18-7 某个值的一个资格迹的例子。访问用星号标记

在 Sarsa 算法中，在时刻 t 的时间误差为：

$$\delta_t = r_{t+1} + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t) \quad (18.18)$$

460 在具有资格迹的 Sarsa 中，称之为 Sarsa(λ)，所有的状态-动作对按下式更新：

$$Q(s, a) \leftarrow Q(s, a) + \eta \delta_t e_t(s, a), \forall s, a \quad (18.19)$$

上式对所有状态-动作对的资格进行更新，更新依赖于其过去出现有多久。 λ 值定义了时间信度：如果 $\lambda = 0$ ，则只进行单步更新。我们在 18.5.3 节讨论的算法就是属于这类，也正因为如此，它们被命名为 $Q(0)$ 、Sarsa(0) 或 TD(0)。随着 λ 趋近于 1，之前的更多步骤被更新。当 $\lambda = 1$ 的时候，所有之前的步骤均被更新，并且分配给它们的信度仅以每步 γ 进行下降。在在线更新中，所有的资格值在每步之后立即更新；而在离线更新中，更新累积至片段结束进行单步更新。在线更新花费更多的时间但是收敛得更快。Sarsa(λ) 的伪代码如图 18-8 所示。 $Q(\lambda)$ 和 TD(λ) 算法可类似的得到 (Sutton 和 Barto 1998)。

```

任意初始化所有的  $Q(s, a)$ ,  $e(s, a) \leftarrow 0, \forall s, a$ 
For 所有的片段
    初始化  $s$ 
    使用由  $Q$  导出的策略 (例如  $\epsilon$ -贪心) 选择  $a$ 
    Repeat
        执行动作  $a$  观测  $r$  和  $s'$ 
        使用由  $Q$  导出的策略 (例如  $\epsilon$ -贪心) 选择  $a'$ 
         $\delta \leftarrow r + \gamma Q(s', a') - Q(s, a)$ 
         $e(s, a) \leftarrow 1$ 
        For 所有的  $s, a$ :
             $Q(s, a) \leftarrow Q(s, a) + \eta \delta e(s, a)$ 
             $e(s, a) \leftarrow \gamma \lambda e(s, a)$ 
             $s \leftarrow s', a \leftarrow a'$ 
    Until  $s$  是终止状态
    
```

图 18-8 Sarsa(λ) 算法

18.6 推广

迄今为止,我们假定 $Q(s, a)$ 值(或者 $V(s)$ 值,如果估计的是状态值)存储在一个查找表中,而我们之前考虑的各种算法称为表格(tabular)算法。这种方法具有几个问题:1)当状态个数和动作个数很大时,表格的尺寸会变得非常大;2)状态和动作可能是连续的,例如,将方向盘以某个角度进行调整;而使用表格,将对这些连续值进行离散化,这可能会导致误差;3)当搜索空间比较大的时候,可能需要非常多的片段才能以可接受的准确程度填满表格的所有项。

我们可以将上述问题看作是一个回归问题,来取代使用表格存储 Q 值。这是一个监督学习问题,其中我们定义一个回归器 $Q(s, a | \theta)$,将 s 和 a 作为输入并通过参数向量 θ 进行参数化来学习 Q 值。例如,这个回归器可以是一个人工神经网络,以 s 和 a 为输入,一个输出,并以 θ 为连接权重。

一个好的函数逼近器具有通常意义上的优势并可以解决之前讨论过的如下问题:一个好的逼近可以用一个简单模型来实现,而不必显式存储训练实例;可以使用连续输入;可以推广:如果我们知道相似的 (s, a) 对之间具有相似的 Q 值,则我们能够对之前的情况进行推广并产生好的 $Q(s, a)$ 值,即使这一状态之前从未遇到过。

为了可以对回归器进行训练,我们需要一个训练集。在Sarsa(0)的情况下,之前我们看到,我们希望 $Q(s_t, a_t)$ 的值最好接近 $r_{t+1} + \gamma Q(s_{t+1}, a_{t+1})$ 的值。这样,我们可以形成一个训练集,其中输入是状态-动作对 (s_t, a_t) ,而要求的输出是 $r_{t+1} + \gamma Q(s_{t+1}, a_{t+1})$ 。我们可以将均方误差写为:

$$E^t(\theta) = [r_{t+1} + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t)]^2 \quad (18.20)$$

可以类似地定义 $Q(0)$ 和TD(0)的训练集。对于后者而言,我们学习 $V(s)$,而要求的输出是 $r_{t+1} + \gamma V(s_{t+1})$ 。一旦准备好训练集,我们可以使用任何监督学习算法在训练集上进行学习。

如果我们使用梯度下降方法,像训练神经网络那样,参数向量可更新如下:

$$\Delta \theta = \eta [r_{t+1} + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t)] \nabla_{\theta} Q(s_t, a_t) \quad (18.21)$$

这是单步更新。而在Sarsa(λ),资格迹也被计算在内:

$$\Delta \theta = \eta \delta_t e_t \quad (18.22) \quad 462$$

其中时间差分误差是:

$$\delta_t = r_{t+1} + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t)$$

并且资格参数向量更新如下:

$$e_t = \gamma \lambda e_{t-1} + \nabla_{\theta} Q(s_t, a_t) \quad (18.23)$$

其中 e_0 为零向量。在表格算法的情况下,会为每对状态-动作存储其资格,因为这些即为(存储为表格)参数。而在使用估计子的情况下,资格是和估计子的参数相关联的。我们也注意到这非常类似于用于稳定反向传播的动量法(11.8.1节)。不同之处在于在动量法中记忆的是先前的权重变化,而这里记忆的是先前的梯度向量。根据计算 $Q(s_t, a_t)$ 所使用的模型,比如神经网络,我们将其梯度向量插入式(18.23)。

理论上,任何回归方法都可用于训练 Q 函数,但是针对这一特定任务还是有若干要求。首先,使用的方法应可以推广,也就是说,我们的确需要保证相似的状态和动作具有相似的 Q 值。同时像在其他应用中一样,也需要对 s 和 a 有一个好的表示,使得相似性比较明显。

其次, 增强学习更新以一个接一个的方式提供实例, 而不是作为一个整体的训练集, 因而学习算法应当有能力进行单个更新来对新的实例进行学习并且不会忘记以前已经学到的东西。例如, 只要使用一个很小的学习率, 一个后向传播多层感知器可通过一个单独的实例进行训练。或者, 可以收集这些实例形成一个训练集来进行学习, 但是这种方法减慢了学习速度, 因为在一个足够大的样本集被收集到之前不会进行任何学习。

由于这些原因, 使用局部学习器对 Q 值进行学习似乎是一个好主意。在这类方法中, 例如径向基函数, 信息被局部化并且当对一个新的实例进行学习的时候, 学习器的一个局部被更新, 而不损坏其他部分的信息。相同的要求也适用于用 $V(s_t | \theta)$ 估计状态值。

463

18.7 部分可观测状态

18.7.1 场景

在某些应用中, 智能主体并不确切的知道系统状态。智能主体配备以传感器, 传感器返回观测 (observation), 而智能主体使用这些观测对系统状态进行估计。比如我们有一个在房间内导航的机器人。这个机器人也许并不知道其在房间内的确切位置, 或还有其他什么东西在房间内。机器人可能装备了一个照相机, 使用它来记录传感观测。虽然这样并不能告诉机器人其确切的状态但是可以提供关于其可能状态的提示信息。例如, 这个机器人可能只知道其右边有一堵墙。

这一场景类似于一个马尔可夫决策过程, 不同之处是在执行动作 a_t 之后, 新的状态 s_{t+1} 是未知的, 但是有一个观察 o_{t+1} , 它是一个关于 s_t 和 a_t 的随机函数: $p(o_{t+1} | s_t, a_t)$ 。这称为部分可观测马尔可夫决策过程 (partially observable MDP, POMDP)。如果 $o_{t+1} = s_{t+1}$, 则 POMDP 简化为 MDP。这就像可观测的和隐马尔可夫模型之间的差别, 而且它们的求解也类似; 也就是说, 我们需要从观测来推断状态 (或状态的概率分布) 并据此执行动作。如果智能主体认为其处于状态 s_1 的概率为 0.4 而处于状态 s_2 的概率为 0.6, 则任一动作的值就是 0.4 乘以在 s_1 状态执行该动作的值加上 0.6 乘以在 s_2 状态执行该动作的值。

对于观测而言, 马尔可夫性质并不成立。下一状态的观测并不仅仅依赖于当前的动作和观测。当只存在有限的观测的时候, 两个状态表面上看起来可能是一样的, 但是实际上却是不同的, 而且如果这两个状态要求执行不同的动作, 那么就会导致以累积奖励为度量的性能上的损失。智能主体应当以某种方式将过去的轨迹压缩到一个当前的单一状态估计。这些过去的观测还可以通过将观测上的一个过去的窗口作为策略输入而被计算在内, 或者使用递归神经网络 (见 11.12.2 节), 在不忘记过去的观测的同时维持状态估计。

在任何时候, 智能主体都可以对最可能的状态进行计算并执行相应的动作。或者它可以执行动作来收集信息并减小不确定性, 例如搜索一个地标, 或停下来询问方向等。这意味着信息价值 (value of information) 的重要性, 并且事实上 POMDP 可以建模为动态 (dynamic) 影响图 (见 16.8 节)。智能主体根据动作所提供的信息、所产生的奖励大小以及它们如何改变环境状态来选择动作。

464

为了保持整个过程是马尔可夫的, 智能主体维护一个内部的信任状态 (belief state) b_t 来对其经历进行总结 (见图 18-9)。智能主体有一个状态估计子 (state estimator), 它基于上一

动作 a_t 、当前观测 o_{t+1} 和前一信任状态 b_t 来更新信任状态 b_{t+1} 。智能主体还有一个策略 π ，与完全可观测环境中所使用的真实状态相反，策略 π 基于这个信任状态来产生下一动作 a_{t+1} 。信任状态是给定初始信任状态(在执行任何动作之前)的环境状态和智能主体以往的观测-动作历史(没有遗漏任何可能提高智能主体性能的信息)上的概率分布。在这种情况下，Q学习使用的是信任状态-动作对的值，而非实际的状态-动作对的值：

$$Q(b_t, a_t) = E[r_{t+1}] + \gamma \sum_{b_{t+1}} P(b_{t+1} | b_t, a_t) V(b_{t+1}) \quad (18.24)$$

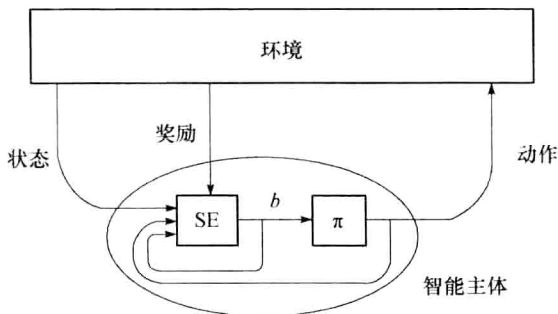


图 18-9 在部分可观测环境中，智能主体具有一个状态估计子(SE)对内部信任状态 b 进行维护并且策略 π 根据这些信任状态产生动作

18.7.2 例子：老虎问题

现在，我们讨论一个例子，它与 Kaelbling、Littmann 和 Cassandra 1998 讨论的，并在 Thrun、Burgard 和 Fox 2005 中修改作为例子的老虎问题(tiger problem)稍微有些不同。假设我们正站在两扇门前，一扇门在我们的左边，而另一扇在右边，分别通往两个房间，两扇门之一的后面潜伏着一只老虎，而另一扇门后有一笔财富。对此，我们并不知晓。如果我们打开有老虎的门，则我们得到很大的负奖励，而如果我们打开有财富的门，则我们得到一些正奖励。隐藏状态 z_t 是老虎的位置，假设 p 表示老虎在左边房间的概率，因而老虎在右边房间的概率为 $1-p$ ：

$$p \equiv P(z_L = 1)$$

两个动作是 a_L 和 a_R ，分别对应于打开左边和右边的门。奖励是

$r(A, Z)$	老虎在左边	老虎在右边
打开左边的门	-100	+80
打开右边的门	+90	-100

我们可以计算这两个动作的期望奖励。不存在进一步的奖励，因为一旦我们打开两扇门中的一扇，场景就结束了。

$$R(a_L) = r(a_L, z_L)P(z_L) + r(a_L, z_R)P(z_R) = -100p + 80(1-p)$$

$$R(a_R) = r(a_R, z_L)P(z_L) + r(a_R, z_R)P(z_R) = 90p - 100(1-p)$$

给定这些奖励，如果 p 接近于 1，即如果我们相信老虎在左边的可能性大，则正确的动作是选择右边的门。类似地，如果 p 接近于 0，则最好选择左边的门。

当 p 在 0.5 附近时，二者交叉，并且期望奖励大约为 -10。事实上，当 p 在 0.5 附近时

(当我们不确定时)期望奖励为负,表明了收集信息的重要性。如果我们能够增加探测设备来降低不确定性,即如果我们能够把 p 从 0.5 移动到 0 或 1 附近,则我们可以采取具有较高正奖励的行动。检测动作 a_s 可以具有较小的负奖励: $R(a_s) = -1$; 这可以看作检测的开销,或看作

466

等价于对未来奖励按 $\gamma < 1$ 打折扣,因为我们推迟了采取(打开一扇门)实际行动。

在这种情况下,期望的奖励和最佳动作的值显示在图 18-10(a)中:

$$V = \max(a_L, a_R, a_s)$$

假设作为检测的输入,我们使用麦克风检测老虎是在左边还是在右边门后,但是我们只有一个不可靠的传感器(因此我们仍然处在部分可观测状态)。假设我们只能以 0.7 的概率检测老虎的存在:

$$\begin{aligned} P(o_L | z_L) &= 0.7 & P(o_L | z_R) &= 0.3 \\ P(o_R | z_L) &= 0.3 & P(o_R | z_R) &= 0.7 \end{aligned}$$

如果我们检测了 o_L ,则我们对老虎位置的看法改变:

$$p' = P(z_L | o_L) = \frac{P(o_L | z_L)P(z_L)}{p(o_L)} = \frac{0.7p}{0.7p + 0.3(1-p)}$$

其效果显示在图 18-10(b)中,在图中我们绘制了 $R(a_L | o_L)$ 。检测到 o_L 使得打开右边的门在更大的范围内成为较好的动作。我们拥有的传感器越好(如果正确检测的概率从 0.7 移近 1),这个范围就越大(习题 9)。类似地,正如我们在图 18-10(c)中看到的,如果我们检测到 o_R ,则提高了打开左边门的可能性。注意,在需要(多次)检测的地方,检测也可能缩小这个范围。

在这种情况下,动作的期望奖励是

$$\begin{aligned} R(a_L | o_L) &= r(a_L, z_L)P(z_L | o_L) + r(a_L, z_R)P(z_R | o_L) \\ &= -100p' + 80(1-p') \\ &= -100 \frac{0.7p}{p(o_L)} + 80 \frac{0.3(1-p)}{p(o_L)} \\ R(a_R | o_L) &= r(a_R, z_L)P(z_L | o_L) + r(a_R, z_R)P(z_R | o_L) \\ &= 90p' - 100(1-p') \\ &= 90 \frac{0.7p}{p(o_L)} - 100 \frac{0.3(1-p)}{p(o_L)} \\ R(a_s | o_L) &= -1 \end{aligned}$$

这种情况下的最佳动作是最大化这三者的动作。类似地,如果我们检测到 o_R ,则期望奖励变成

$$\begin{aligned} R(a_L | o_R) &= r(a_L, z_L)P(z_L | o_R) + r(a_L, z_R)P(z_R | o_R) \\ &= -100 \frac{0.3p}{p(o_R)} + 80 \frac{0.7(1-p)}{p(o_R)} \\ R(a_R | o_L) &= r(a_R, z_L)P(z_L | o_R) + r(a_R, z_R)P(z_R | o_R) \\ &= 90 \frac{0.3p}{p(o_R)} - 100 \frac{0.7(1-p)}{p(o_R)} \\ R(a_s | o_R) &= -1 \end{aligned}$$

467

为了计算期望奖励,我们需要在两个传感器读数上用它们的概率取加权平均:

$$\begin{aligned}
 V' &= \sum_j [\max_i R(a_i | o_j)] P(o_j) \\
 &= \max(R(a_L | o_L), R(a_R | o_L), R(a_S | o_L)) P(o_L) + \\
 &\quad \max(R(a_L | o_R), R(a_R | o_R), R(a_S | o_R)) P(o_R) \\
 &= \max(-70p + 24(1-p), 63p - 30(1-p), -0.7p - 0.3(1-p)) + \\
 &\quad \max(-30p + 56(1-p), 27p - 70(1-p), -0.3p - 0.7(1-p)) \\
 &= \max \begin{pmatrix} -100p & +80(1-p) \\ -43p & -46(1-p) \\ 33p & +26(1-p) \\ 90p & -100(1-p) \end{pmatrix} \quad (18.25)
 \end{aligned}$$

468

注意, 当我们乘以 $P(o_L)$ 时, 它被约去, 而我们得到 p 的线性函数。这五条线和对应于它们的最大值的分段函数显示在图 18-10(d) 中。注意直线 $-40p - 5(1-p)$, 以及涉及 a_S 的直线都在其他直线下方, 因此可以安全地删除。图 18-10(d) 比图 18-10(a) 更好这一事实表明了信息的价值。

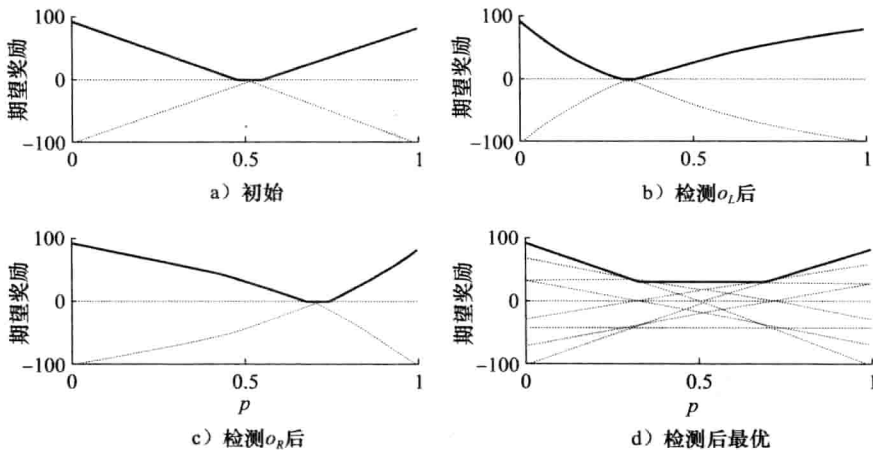


图 18-10 老虎问题的期望奖励和检测的效果

这里, 我们要计算的是选取 a_S 后的最佳动作值。例如, 第一条线对应于选取 a_S 之后选择 a_L 。因此, 为了找到长度为 2 的片断的最佳决策, 我们需要通过减 1 来实现 (-1 是 a_S 的奖励), 并得到检测动作的期望奖励。等价地, 我们可以把这看作具有立即奖励 0 但将未来的奖励按 $\gamma < 1$ 打折的等待。我们还有两个通常的动作 a_L 和 a_R , 并且我们选择这三个中的最佳的; 两个立即动作和一个打折的未来动作。

现在, 让我们像 Thrun、Burgard 和 Fox 2005 的例子那样, 将问题设计得更有趣。我们假设两个房间之间有一扇门, 而我们看不到, 老虎可以从一个房间移动到另一个。假设这是一只好动的老虎, 它待在一个房间的概率为 0.2, 而走到另一个房间的概率为 0.8。这意味着 p 也应该更新为

$$p' = 0.2p + 0.8(1-p)$$

并且在选择 a_S 之后选择最佳动作时, 在式 (18.25) 中使用这个更新后的 p :

$$V' = \max \begin{pmatrix} -100p & +80(1-p) \\ 33p & +26(1-p) \\ 90p & -100(1-p) \end{pmatrix}$$

469 图 18-11(b) 对应于使用更新后 p' 的图 18-10(d)。现在, 当我们规划长度为 2 的片断时, 有两个立即动作 a_L 和 a_R , 或者当 p 改变时我们等待并检测, 然后执行动作并得到它的打折的奖励(参见图 18-11b):

$$V_2 = \max \begin{pmatrix} -100p & +80(1-p) \\ 90p & -100(1-p) \\ \max V' - 1 \end{pmatrix}$$

我们看到图 18-11(b) 比图 18-10(a) 好: 当错误动作可能导致大的惩罚时, 最好推迟决断, 寻求附加信息, 并提前规划。我们可以通过以下方法来考虑更长的片断: 继续 p 的迭代更新, 通过减 1 打折, 并包含两个立即动作 a_L 和 a_R , 以计算 $V_t, t > 2$ 。

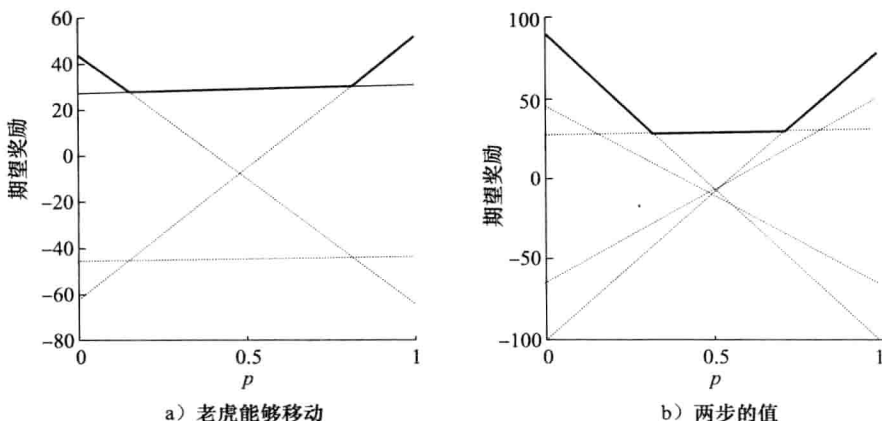


图 18-11 期望奖励改变: a) 如果隐藏状态可以改变; b) 考虑长度为 2 的片断

我们刚讨论的值用分段线性函数表示的算法只能在状态数、动作数、观测数和片断长度均为有限时才能使用。即便在这些值的某个不太小, 或存在连续值的应用中, 算法的复杂度都很高, 因而我们需要借助于具有合理复杂度的近似算法, 这类算法的综述在 Hauskrecht 2000 以及 Thrun、Burgard 和 Fox 2005 中。

18.8 注释

Thrun、Burgard 和 Fox 2005 包含增强学习用于机器人的近期工作, 还给出了一些令人印象深刻的应用。关于增强学习的更多信息可以在 Sutton 和 Barto (1998) 的教科书中找到, 该书讨论了增强学习的各个方面、学习算法以及若干应用。而 Kaelbling、Littman 和 Moore 1996 是增强学习的全面介绍。

Bertsekas 1987 以及 Bertsekas 和 Tsitsiklis 1996 讨论了动态规划方法, 而 Q 学习可以看作是动态规划的随机近似 (Jaakkola、Jordan 和 Singh 1994)。与经典的动态规划相比, 增强学习有两个优点: 首先, 在学习期间, 增强学习可专注于空间的重要部分而忽略其他部分; 其次, 增强学习可以使用函数逼近方法来表示知识, 进而得以推广和更快地学习。

一个相关的领域是学习自动机(learning automata)(Narendra 和 Thathachar 1974)它是一个有限状态机器,通过“试错”解决类似于 K -臂赌博机问题。我们这里所讨论的场景同样也是最优控制的课题,其中一个控制器(智能主体)在设施(环境)中执行动作来最小化系统开销(最大化奖励)。

最早使用时间差分方法的是 Samuel 写于 1959 年的跳棋游戏程序(Sutton 和 Barto 1998)。对于一个游戏中每对相继的位置,通过棋盘评估函数对两个棋盘状态进行评估,进而引发一个更新来减小它们之间的差异。关于游戏方面的研究工作很多,因为其兼具易于定义和挑战性的特点。对一个类似象棋的游戏的模拟也易于进行:允许的棋步可以形式化而且目标状态清晰。尽管定义这样一个游戏很简单,但是以专家级别进行游戏却非常困难。

增强学习最令人印象深刻的应用是 TD-Gammon 程序。该程序通过和自身进行对弈来学习西洋双陆棋(Tesauro 1995)。它优于同样由 Tesauro 开发的 neruogammon 程序,后者基于与专家对弈,以监督学习方式进行训练。西洋双陆棋是大约有 10^{20} 种状态的复杂任务,并存在由于掷骰子而产生的随机性。使用 TD(λ) 算法,TD-Gammon 程序在经过和自身副本 1 500 000 次对弈后达到了大师级水平。

另一个有趣的应用是作业车间调度(job shop scheduling)问题或寻找满足时间和资源约束的任务调度问题(Zhang 和 Dietterich 1996)。某些任务必须要在其他任务开始之前完成,并且需要相同资源的两个任务不能同时进行。Zhang 和 Dietterich 使用增强学习很快找到了满足约束并且较短的调度方式。每个状态是一个调度,而动作是调度更改,最终程序找到的不仅是一个好的调度,而且是对一类相关调度均有效的调度。

471

最近提出了层次化方法将问题分解为一组子问题。其优点是针对子问题学习而得到的策略可在多个问题上共享,这加速了对新问题的学习速度(Dietterich 2000)。每个子问题都更简单,并且对它们单独进行学习更快一些。缺点是当对子问题的策略进行组合时,所得的策略可能是次最优的。

尽管增强学习算法比监督学习算法慢一些,但很明显它们具有更广泛的应用并具有构建更好学习机器的潜力(Ballard 1997)。它们不需要任何监督,因而可能实际上更好一些,因为不会被老师误导。例如, Tesauro 的 TD-Gammon 程序在某些情况下所走的棋步比最好的棋手所走的棋步还要好。增强学习领域发展迅速,因而我们可以期待在不远的将来看到其他引人注目的成果。

18.9 习题

1. 给定图 18-12 的网格世界,如果到达目标的奖励为 100 并且 $\gamma = 0.9$,手工计算 $Q^*(s, a)$, $V^*(S)$ 以及最优策略的动作。

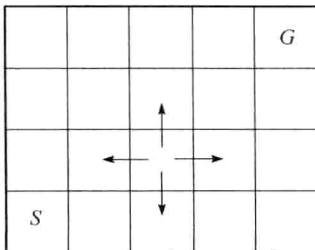


图 18-12 网格世界。智能主体始于 S , 可以向四个罗盘方向移动。目标状态为 G

- 472 2. 以习题 1 中相同的配置, 使用 Q 学习算法学习最优策略。
3. 在习题 1 中, 如果在右下角加入另一个目标状态, 最优策略将如何改变? 如果在右下角的状态定义奖励为 -100 (非常坏的状态) 将发生什么?
4. 作为对 $\gamma < 1$ 的替代, 有 $\gamma = 1$ 并且所有中间状态 (非目标) 具有一个负的奖励 $-c$ 。这二者有何差异?
5. 在习题 1 中, 假设到达目标的奖励服从均值 100 和方差 40 的正态分布。同时假设动作也是随机的, 即当机器人向一个方向前进的时候, 它以 0.5 的概率向预定的方向前进同时以 0.25 的概率向两个横向方向之一前进。在这种情况下, 学习 $Q(s, a)$ 。
6. 假设我们想要使用 $TD(\lambda)$ 算法对状态值函数 $V(s)$ 进行估计。推导出其表值迭代更新。
7. 使用式 (18.22), 推导出使用多层感知器估计 Q 的权重更新公式。
8. 给出一个可用 POMDP 建模的增强学习应用的例子。定义其中的状态、动作、观测和奖励。
9. 在老虎例子中, 说明当我们有更可靠的检测设备时, 在需要一再检测的地方, 范围将减小。
10. 使用如下奖励矩阵重做老虎例子:

$r(A, Z)$	老虎在左边	老虎在右边
打开左边的门	-100	+10
打开右边的门	20	-100

18.10 参考文献

- Ballard, D. H. 1997. *An Introduction to Natural Computation*. Cambridge, MA: MIT Press.
- Bellman, R. E. 1957. *Dynamic Programming*. Princeton: Princeton University Press.
- Bertsekas, D. P. 1987. *Dynamic Programming: Deterministic and Stochastic Models*. New York: Prentice Hall.
- Bertsekas, D. P., and J. N. Tsitsiklis. 1996. *Neuro-Dynamic Programming*. Belmont, MA: Athena Scientific.
- 473 Dietterich, T. G. 2000. "Hierarchical Reinforcement Learning with the MAXQ Value Decomposition." *Journal of Artificial Intelligence Research* 13: 227-303.
- Hauskrecht, M. 2000. "Value-Function Approximations for Partially Observable Markov Decision Processes." *Journal of Artificial Intelligence Research* 13: 33-94.
- Jaakkola, T., M. I. Jordan, and S. P. Singh. 1994. "On the Convergence of Stochastic Iterative Dynamic Programming Algorithms." *Neural Computation* 6: 1185-1201.
- Kaelbling, L. P., M. L. Littman, and A. R. Cassandra. 1998. "Planning and Acting in Partially Observable Stochastic Domains." *Artificial Intelligence* 101: 99-134.
- Kaelbling, L. P., M. L. Littman, and A. W. Moore. 1996. "Reinforcement Learning: A Survey." *Journal of Artificial Intelligence Research* 4: 237-285.

- Narendra, K. S., and M. A. L. Thathachar. 1974. "Learning Automata—A Survey." *IEEE Transactions on Systems, Man, and Cybernetics* 4: 323-334.
- Sutton, R. S. 1988. "Learning to Predict by the Method of Temporal Differences." *Machine Learning* 3: 9-44.
- Sutton, R. S., and A. G. Barto. 1998. *Reinforcement Learning: An Introduction*. Cambridge, MA: MIT Press.
- Tesauro, G. 1995. "Temporal Difference Learning and TD-Gammon." *Communications of the ACM* 38(3): 58-68.
- Thrun, S., W. Burgard, and D. Fox. 2005. *Probabilistic Robotics*. Cambridge, MA: MIT Press.
- Watkins, C. J. C. H., and P. Dayan. 1992. "Q-learning." *Machine Learning* 8: 279-292.
- Zhang, W., and T. G. Dietterich. 1996. "High-Performance Job-Shop Scheduling with a Time-Delay TD(λ) Network." In *Advances in Neural Information Processing Systems 8*, ed. D. S. Touretzky, M. G. Mozer, and M. E. Hasselmo, 1024-1030. Cambridge, MA: The MIT Press.

机器学习实验的设计与分析

我们讨论评估和比较实际学习算法性能的机器学习实验的设计，并分析这些实验结果的统计检验。

19.1 引言

在前面的章节中，我们讨论了一些学习算法，并且看到对于一个给定的应用，多种算法都是可行的。现在，我们关心的是以下两个问题：

1) 如何评估一个学习算法在给定问题上的期望误差？也就是说，例如，已经使用分类算法在取自某个应用的数据集上训练了一个分类器，我们是否能够以足够的置信度认为在之后的实际应用中，其期望误差率将小于，比如2%？

2) 给定两个学习算法，就给定的应用而言，我们能够认为一个算法的误差比另一个低吗？进行比较的分类算法可能是不同的，例如，参数的与非参数的，或者它们可能使用不同的超参数设置。例如，给定一个具有4个隐藏单元的多层感知器(参见第11章)和另一个具有8个隐藏单元的感知器，我们希望可以判断哪一个具有更低的期望误差；或者在使用 k -最近邻分类器(参见第8章)时，我们希望找到最佳的 k 值。

我们不能只看训练集上的误差并据此来做判定。根据定义，训练集上的误差率总会小于包含训练时未见过的实例的检验集上的误差率。类似地，训练误差不能用于比较两个算法。[475] 因为在训练集上，具有更多参数的复杂模型几乎总是比简单模型的误差更小。

因此，正如我们反复讨论的，我们需要一个不同于训练集的验证集，并且即使是在验证集上，一轮运行也可能不够。其原因有二：首先，训练集和验证集都可能较小并且可能包含异常实例，如噪声或离群点，可能会误导我们；第二，学习算法有可能依赖于影响泛化的其他随机因素。例如，对于使用后向传播训练的一个多层感知器，由于梯度下降收敛于局部极小，所以初始权重会影响最终权重，并且以完全相同的结构和训练集，以不同的初始权重开始训练，最终有可能产生多种分类器，这些分类器在相同的验证集上有不同的错误率。因而我们需要多轮运行，以便平均这些随机源。如果我们只是训练和验证一次，则无法检验这些因素的影响；只有在学习方法的代价很高以至于只能训练和验证一次时，训练和验证一次才是可以接受的。

我们在一个数据集上运行学习算法(learning algorithm)并产生一个学习器(learner)。如果我们只训练一次，则只得到一个学习器和一个验证误差。为了平均各种随机性(来自训练数据、初始权重等)，我们使用相同的算法来产生多个学习器，进而在多个验证集上检验它们并记录验证误差的一个样本(当然，所有训练和验证集均应取自同一应用)。我们对学习算法的评估基于这些验证误差的分布(distribution)。我们可以使用这一分布来评估学习算法在该问题上的期望误差(expected error)，或者将它与某种其他学习算法的误差率分布进行

比较。

在讨论这一过程如何完成之前，需要重点强调以下几点：

1) 需要牢记的是：无论我们从分析中获得何种结论，该结论都限于所给定的数据集。我们并不用独立于领域的方式来比较学习算法，而是针对某一特定应用进行比较。一般而言，我们不对学习算法的期望误差率做任何讨论，也不将一个学习算法和另一个进行比较。我们得到的任何结果只对特定的应用成立，而且仅在该应用可以由我们使用的样本代表的意义上有效。而且无论如何，正如没有免费的午餐法则 (No Free Lunch Theorem) (Wolpert 1995) 所述，没有诸如“最好的”学习算法之说。对于任何学习算法，均有一个数据集使其非常准确，而另一个数据集使其非常差。当我们认为一个学习算法好时，只是量化其归纳偏倚在多大程度上与数据的性质一致。

476

2) 将给定数据集划分为一定数量的训练集和验证集对仅仅是为了检验。一旦所有的检验完成，并且决定了最终方法或超参数，为了训练最终学习器，我们可以使用先前用于训练或验证的所有已标记数据。

3) 由于我们还使用验证集进行检验，比如，为了选择两个学习算法中较好的一个，或决定何时停止学习，验证集实际上成为我们所使用数据的一部分。在结束所有的检验之后，我们选定了某一特定的算法并且希望报告其期望误差，为此我们应使用另外一个在训练最终系统过程中未曾使用过的检验集 (test set)。该数据应当在之前的训练或验证过程中从未使用过，并且应足够大以使误差估计有意义。因此，给定一个数据集，我们应当保留一部分数据作为检验集，而其余的数据用于训练和验证。通常，像我们在稍后看到的那样，我们可以留 1/3 的样本作为检验集，使用另外 2/3 做交叉验证，以产生多对训练/验证集。因而，给定特定学习算法和模型结构，训练集用于参数优化；验证集用于优化学习算法或模型结构的超参数；而一旦二者均被优化，才在最后使用检验集。例如，对一个多层感知器 (MLP) 而言，训练集用于优化权重，验证集用于确定隐藏单元个数、训练多久、学习率等。一旦选择了最佳的 MLP 配置，其最终的误差在检验集上计算。对于 k -NN，训练集作为查找表存储；我们在验证集上优化距离度量和 k 值，最后在检验集上进行检验。

4) 通常，我们就错误率对学习算法进行比较，但应牢记，在现实中，误差仅仅是影响决策的一个标准。一些其他标准是 (Turney 2000)：

477

- 当使用损失函数，而非 0/1 损失 (参见 3.3 节) 对误差进行泛化时的风险。
- 训练的时间和空间复杂度。
- 检验的时间和空间复杂度。
- 可解释性，即使用的方法是否允许提取可以由专家检查和确认的知识。
- 易于编程。

这些因素相对重要程度依赖于应用。例如，如果在工厂中只进行一次训练，那么训练的时间和空间复杂度就不重要；如果在使用过程中要求自适应性，则训练的时间和空间复杂度就变得重要了。多数学习算法使用 0/1 损失并以误差最小化为唯一标准；最近，提出了这些算法的变种，即代价敏感学习 (cost-sensitive learning) 算法把其他代价标准也考虑在内。

当我们在一个数据集上使用训练集训练学习器，在验证集上检验它的准确率并试图提取结论时，我们所做的是实验。统计学提供了系统的方法，指导我们正确地设计实验，告诉我

们如何分析收集的数据,以便得到统计显著的结论(Montgomery 2005)。本章我们将考察如何将这些方法用于机器学习。

19.2 因素、响应和实验策略

与科学和工程的其他分支一样,在机器学习中我们也做实验,以便获得关于考察的过程信息。这里,所考察的是学习器,它已经在数据集上训练,能够对给定的输入产生输出。实验(experiment)是一次检验或一系列检验,其中包含一些影响输出的因素(factor)。这些因素可以是使用的算法、训练集、输入特征等,而我们观察响应(response)的变化,以便能够提取信息。目标可以是识别最重要的因素,排除不重要的因素,或者是找出某种因素配置来优化响应,例如,优化在给定检验集上的分类准确率。

我们的目标是规划和进行机器学习实验,并分析实验产生的数据,以便能够排除随机性的影响,得到统计显著的(statistically significant)结论。在机器学习中,我们的目标是具有最高泛化准确率、最小复杂度的学习器(以便其实现的时间和空间代价低),并且该学习器是鲁棒的,即受外部变化的影响最小。

一个训练后的学习器可能如图 19-1 所示。对于一个检验输入,它产生一个输出(例如,类编码),并且依赖于两类因素:可控因素和不可控因素。正如名称所示,可控因素(controllable factor)是我们控制的那些因素。最基本的可控因素是使用的学习算法。还有算法的超参数,例如,多层感知器的隐藏单元数、 k 最近邻的 k 、支持向量机的 C 等。所使用的数据集和输入表示(即输入如何表示成向量)都是可控因素。

还有一些我们不能控制的不可控因素(uncontrollable factor),将不期望的可变性添加到过程中,我们不想让它们影响我们的决策。这些因素包括数据中的噪声、重复从大数据集中抽样产生的特定训练子集、优化过程中的随机性,例如多层感知器的梯度下降中的初始状态等。

我们使用输出来产生响应(response)变量;例如,检验集上的平均分类误差,使用损失函数的期望风险或其他测度,如我们稍后将讨论的精度和召回率。

给定多种因素,我们需要为最佳响应找出这些因素的最佳设置;或者更一般地,我们需要确定它们对响应变量的影响。例如,在使用 k 最近邻(k -NN)分类方法之前,我们可能使用主成分分析(PCA)将维度降低到 d 。 d 和 k 是两个因素,而问题是 d 和 k 哪个组合导致最优性能。或者,我们可以使用具有高斯核的支持向量机分类方法,而我们需要将正则化参数 C 和高斯分布的展宽 s^2 同时调整。

存在多种实验策略(strategies of experimentation),如图 19-2 所示。在最佳猜测(best guess)方法中,我们从某个我们相信是良好配置的因素设置开始。在此检验响应,并且每次稍微改动一个(或少量)因素,检验每个组合,直至得到一个我们认为足够好的状态。如果实验过程很直观,这种方法可能有效。但要注意,这里没有系统的方法来修改因素,并且当

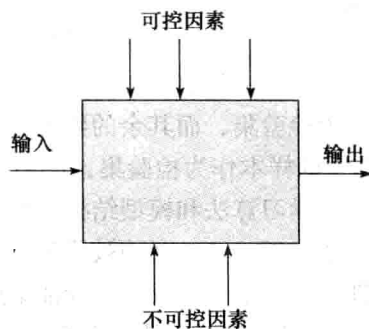


图 19-1 给定输入,产生输出的过程,它受可控因素和不可控因素的影响

我们停止时，不能保证找到最佳配置。

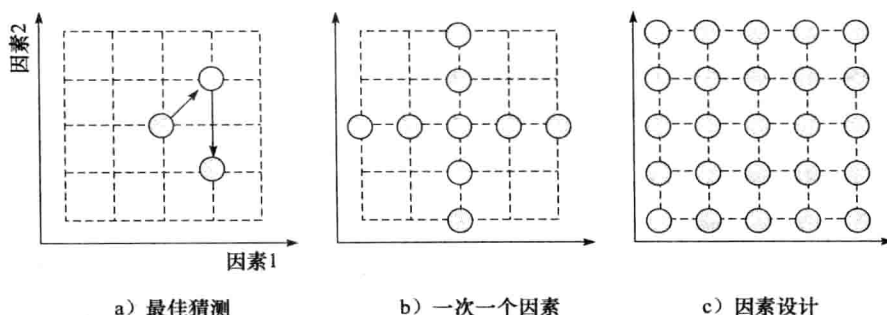


图 19-2 两个因素，每个五个不同水平的实验策略

另一种策略是一次一个因素(one factor at a time)地修改。在这种方法中，我们为所有因素确定一个基线(缺省)值，然后对一个因素尝试不同水平，同时令其他因素保持在基线上。这种方法的主要缺点是，它假定因素之间不相互影响，这一假定并非总能成立。在稍早讨论的 PCA/ k -NN 级联中， d 的每个选择都为 k -NN 定义了不同的输入空间，其中不同的 k 值可能更合适。

480

正确的方法是使用因素设计(factorial design)，其中因素一起变化，不是一次一个；这通俗地称作网格搜索(grid search)。对于 F 个因素，每个因素有 L 个水平，一次一个因素搜索需要 $\mathcal{O}(L \cdot F)$ 时间，而因素设计实验需要做 $\mathcal{O}(L^F)$ 时间。

19.3 响应面设计

为了减少需要的运行次数，一种可能的方法是运行部分因素设计，仅运行所有配置的一个子集；另一种方法是尝试使用上一次运行收集的知识来估计看上去会有高响应的配置。在一次一个因素的搜索中，如果我们能够假定响应通常是二次的(具有单个最大，假定我们最大化响应，如检验准确率)，则不尝试所有的值，而是使用一个迭代过程，从某个初始运行开始，拟合一个二次曲面，解析地找出它的极大值，取它作为下一个估计，在这种情况下运行一次实验，把结果数据添加到样本中，然后继续拟合和抽样，直到不能进一步改进为止。

使用多个因素，这被推广为响应面设计(response surface design)方法，那里我们尝试用有参的响应函数拟合这些因素：

$$r = g(f_1, f_2, \dots, f_F | \phi)$$

其中， r 是响应，而 $f_i (i=1, \dots, F)$ 是因素。这个由给定参数 ϕ 定义的被拟合的有参函数是我们的经验模型，它对(可控的)因素的具体配置估计响应；不可控因素的影响按噪声建模。 $g(\cdot)$ 是一个回归模型(通常是二次的)，并且在基线(由所谓的设计矩阵(design matrix)定义)附近的少数几次运行之后，我们就能得到足够的样本来拟合 $g(\cdot)$ 。然后，我们就能解析地计算使拟合的 g 取最大值的 f_i ，这被我们取作下一个猜测，在下一个猜测中运行实验，得到数据实例，将它添加到样本中，再次拟合 g ，如此下去，直到收敛。这种方法是否有效取决于响应是否确实是因素的具有单个最大值的二次函数。

481

19.4 随机化、重复和阻止

现在，我们讨论实验设计的三项基本原则：

- 随机化(randomization)要求实验运行的顺序应该随机确定，从而使得结果是独立的。通常，在涉及物理对象的现实世界的实验中，这是一个问题。例如，机器需要一些时间来预热直到它们在正常范围运转，因此检验应该在随机次序下完成，以不影响结果。在软件实验中，顺序一般不是问题。
- 重复(replication)意味着应该对(可控的)因素的相同配置做多次实验，以便平均不可控因素的影响。在机器学习中，这通常通过在相同数据集的许多再抽样的版本上运行相同的算法来实现；这称作交叉验证(cross-validation)，我们将在 19.6 节讨论。响应在相同实验的这些不同重复上的变化，使得我们可以得到实验误差(不可控因素的影响)的估计，这又可以用来确定多大的差别才能视为统计显著的(statistically significant)。
- 阻止(blocking)用来降低或消除有害因素(nuisance factor)导致的可变性，有害因素是影响响应但我们对其不感兴趣的因素。例如，工厂的生产缺陷也可能与原材料的批次有关，而这种影响应该与工厂里的设备、人员等可控因素分开。在机器学习实验中，当我们使用再抽样，对不同的重复实验使用数据的不同子集时，需要确保如果我们比较学习算法，则它们应该使用相同的再抽样子集的集合，否则准确率的差异不仅取决于不同的算法，而且还取决于不同的子集。为了能够度量仅由于算法导致的差别，重复运行的不同训练集应该是相同的，这就是阻止的含义。在统计学中，如果有两个总体，则这称作配对(pairing)并用于配对检验(pairing testing)。

482

19.5 机器学习实验指南

在开始实验之前，我们需要清楚研究什么，如何收集数据，我们打算怎样分析它。对任何类型的实验来说，机器学习的步骤都是相同的(Montgomery 2005)。注意，在这里，任务是分类还是回归，是非监督的还是增强学习的应用并不重要。整个讨论都适用，唯一的区别是所收集的响应数据的抽样分布。

1) 研究目标

我们需要通过清楚地陈述问题、定义研究目标开始。在机器学习中，可能存在多种可能性。例如，正如我们以前所讨论的，我们感兴趣的可能是评估一种学习算法在特定问题上的期望误差(或某种其他响应度量)，并检查该误差是否低于某个可以接受的水平。

给定两个学习算法和一个由数据集定义的具体问题，我们可能希望确定哪个算法具有较低的泛化误差。这两个算法可能是不同的算法，也可能一个算法是另一个的改进，例如，通过使用更好的特征提取来改进。

在一般情况下，可能有多个学习算法，而我们可能想要选择具有最低误差的算法，或者对于给定数据集，将它们按误差排序。

在更一般的情况下，我们可能希望在两个或多个数据集上，而不是在单个数据集上比较两个或多个算法。

2) 响应变量的选择

我们需要确定我们应该使用什么作为质量度量。最常使用的是误差，即分类的误分类错误和回归的均方误差。我们也可以使用某些变形，例如，可以使用风险度量，将 0/1 损失推广为任意损失。在信息检索中，我们使用诸如精度和召回率这样的度量，我们将在 19.7 节讨论。在代价敏感学习中，不仅要考虑输出，还要考虑系统参数，例如复杂度。

483

3) 因素和水平的选择

因素是什么取决于研究目标。如果我们固定算法并且想要找出最佳超参数，则这些超参数就是因素。如果我们想比较算法，则学习算法是因素。如果我们有不同的数据集，则这些数据集也成为因素。

因素的水平应该小心选择，以便不失去好的配置，并且避免做不必要的实验。最好试着对因素水平规范化。例如，在优化 k 最近邻的 k 值时，可以尝试 1、3、5 等值，但在优化 Parzen 窗口的展宽 h 时，我们不要尝试诸如 1.0、2.0 等绝对值，因为它依赖于输入的标度。最好使用指示标度的统计量，例如，实例与它的最近邻之间的平均距离，并且尝试将 h 作为该统计量的不同倍数。

尽管以前的经验一般是加分的，但是同样重要的是考察所有的因素和可能重要的因素水平，而不过于受以往经验的影响。

4) 实验设计的选择

除非我们确信因素之间不相互影响，否则最好做因素设计，因为因素之间多半会相互影响。重复的次数依赖于数据集的规模；当数据集大时，它可以保持小；下一节讨论再抽样时我们将讨论这一问题。然而，太少的重复产生少量数据，使得分布比较很困难；在参数检验这种特殊情况下，高斯分布假设可能靠不住。

一般地，给定一个数据集，我们留一部分作为检验集，而其余的用来训练和验证，在大部分时候这可以通过再抽样来做。如何进行划分是重要的，实践中，使用小数据集导致具有高方差的响应，差别可能是不显著的，并且结果可能不是令人信服的。

此外，重要的是尽可能地避免使用无使用价值的人造数据集，而要使用从现实生活环境收集的的实际数据集。教学用的一、二维数据集可能有助于提供直观解释，但是在高维空间，算法的行为可能完全不同。

484

5) 做实验

在运行具有许多因素和水平的大型因素实验之前，最好先对某些随机设置试运行几次，检查一切是否如预期的一样。在一个大型实验中，最好保留一些中间结果（或随机数产生器的种子），以便需要时可以重新运行整个实验的一部分。所有的结果都应当是可再现的。在运行具有许多因素和因素水平的大型实验时，应该清楚软件老化的负面影响。

重要的是，在做实验时，实验者是无偏向的。在将个人喜爱的算法与其他算法比较时，两个算法都要同样仔细地考察。在大规模的研究中，甚至可以设想测试者不同于开发者。

应该避免写自己的“程序库”的诱惑，而是应该尽可能地使用来源可靠的程序，这样的程序经受了更好的测试和优化。

与任何软件开发研究一样，好文档的作用不可低估，特别是在分组开发时。高质量的软件工程开发的所有方法也应该用在机器学习实验中。

6) 数据的统计分析

这对应于用这样的方式分析数据：无论我们得到什么结论，这些结论都不是主观的或随机的。我们在假设检验的框架下提出想要回答的问题，并检查样本是否支持该假设。例如，问题“ A 是比 B 更准确的算法吗？”变成假设“我们能够说 A 训练的学习器的平均误差显著低于 B 训练的学习器吗？”

像通常一样，可视分析是有帮助的，并且我们可以使用误差分布的直方图、盒图、变程图 485 图等。

7) 结论和建议

一旦收集和分析了所有数据，我们就可以提取客观的结论。最常遇到的结论是需要进一步实验。大部分统计学研究、大部分机器学习和数据挖掘研究都是迭代的。正因为如此，我们从来都不是一开始就做所有的实验。有人建议，第一次实验考察的数据不超过可利用资源的 25% (Montgomery 2005)。第一次实验只是调查。这就是为什么开始最好不要抱太大期望，向你的老板或论文导师承诺什么的原因。

我们应该始终牢记，统计检验不会告诉我们假设是否正确，而是指出样本看上去与假设的一致程度有多大。总是存在得不到结论性结果或者得到错误结论的风险，特别是当数据集很小和存在噪声时。

当我们的期望不满足时，最好是考察为什么它们不满足。例如，检查为什么我们钟爱的算法 A 在某些情况下效果极差，我们可能对 A 改进版本产生绝妙的想法。所有的改进都是由于以前的版本有缺陷；找到缺陷是有益的暗示：存在我们可以做的改进！

但是，在我们确信已经完全分析了当前的数据并且从中学习到我们能够学习的一切之前，不要急于做改进版本的下一步检验。想法是廉价的、无用的，除非它被检验，而检验是昂贵的。

19.6 交叉验证和再抽样方法

为了重复实验，我们的第一个需求是(在留下一些作为检验集后)从数据集 \mathcal{X} 中获得一定数目的训练集和验证集对。为此，如果样本 \mathcal{X} 足够大，我们可以随机地将其分为 K 个部分，然后将每一部分随机地分为两部分，一半用于训练，另一半用于验证。 K 通常为 10 或 30。不幸的是，数据集从未有如此之大，以允许我们这样做。因此，我们应该在小数据集上尽力而为，其方法是以不同划分来重复使用相同数据，这称为交叉验证(cross-validation)。其潜 486 在的问题是交叉验证使得错误率是相互依赖的，因为这些不同集合共享了数据。

因此，给定一个数据集 \mathcal{X} ，我们希望通过该数据集产生 K 对训练和验证集 $\{\mathcal{T}_i, \mathcal{V}_i\}_{i=1}^K$ 。我们要保持训练和验证集尽可能大，以保证误差估计的鲁棒性；同时，要保持不同集合间的重叠尽可能小。我们还要确保当抽取数据子集的时候，类以正确比例被代表，不扰乱类的先验概率，这称为分层(stratification)。如果一个类在整个数据集中占有 20% 的实例，则在所有取自该数据集的抽样集中，该类也应该大约有 20% 的实例。

19.6.1 K -折交叉验证

在 K -折交叉验证(K -fold cross-validation)中，数据集 \mathcal{X} 被随机地划分为 K 等份 $\mathcal{X}_i, i=1,$

..., K 。为了产生每对训练集和验证集,我们将 K 份数据中的一份保留为验证集,其余 $K-1$ 份合并为训练集。重复 K 次,每次保留 K 份中的另一份数据,可得到 K 对数据集:

$$\begin{aligned}\mathcal{V}_1 &= \mathcal{X}_1 & \mathcal{T}_1 &= \mathcal{X}_2 \cup \mathcal{X}_3 \cup \cdots \cup \mathcal{X}_K \\ \mathcal{V}_2 &= \mathcal{X}_2 & \mathcal{T}_2 &= \mathcal{X}_1 \cup \mathcal{X}_3 \cup \cdots \cup \mathcal{X}_K \\ &\vdots & & \\ \mathcal{V}_K &= \mathcal{X}_K & \mathcal{T}_K &= \mathcal{X}_1 \cup \mathcal{X}_2 \cup \cdots \cup \mathcal{X}_{K-1}\end{aligned}$$

这种方法有两个问题:首先,为了保持训练集较大,我们允许验证集较小;其次,训练集在相当大的程度上重叠,确切地说,任意两个训练集共享 $K-2$ 份数据。

K 一般为 10 或 30。当 K 增加的时候,用于训练实例的比例增加因而产生更为鲁棒的估计,但是验证集相应变小。此外,也带来了将分类器训练 K 次的代价,这一代价随着 K 增加而增加。当 N 增大时, K 可以较小;如果 N 小,则 K 应该大,以保证有足够大的训练集。 K -折交叉验证的一个极端情况是留一(leave-one-out),其中给定 N 个实例的数据集,只保留一个实例作为验证集(验证实例),其余 $N-1$ 个实例作为训练集。由此我们通过每次迭代中保留一个不同的实例而得到 N 对不同的训练集和验证集。这种方法通常用于诸如医疗诊断的应用中,这类应用很难找到标记数据。留一无法保证分层。

最近,随着计算费用的降低,多次运行 K -折交叉验证已经成为可能(例如, 10×10 折),并且在平均值上取平均,以便得到更可靠的误差估计(Bouchaert 2003)。

19.6.2 5×2 交叉验证

Dietterich(1998)提出了 5×2 交叉验证(5×2 cross-validation),使用等大小的训练集和验证集。将数据集 \mathcal{X} 随机地分为两部分: $\mathcal{X}_1^{(1)}$ 和 $\mathcal{X}_1^{(2)}$,这样就给出了第一对训练集和验证集: $\mathcal{T}_1 = \mathcal{X}_1^{(1)}$ 和 $\mathcal{V}_1 = \mathcal{X}_1^{(2)}$ 。然后我们交换两个半份的角色来得到第二对训练集和验证集: $\mathcal{T}_2 = \mathcal{X}_1^{(2)}$ 和 $\mathcal{V}_2 = \mathcal{X}_1^{(1)}$ 。这就是第一次对折; $\mathcal{X}_i^{(j)}$ 表示第 i 次对折中的第 j 个半份。

为了得到第二次对折,我们随机地将 \mathcal{X} 打乱并将其划分为新的对折 $\mathcal{X}_2^{(1)}$ 和 $\mathcal{X}_2^{(2)}$ 。这可通过从 \mathcal{X} 中随机无放回抽样来实现,即 $\mathcal{X}_1^{(1)} \cup \mathcal{X}_1^{(2)} = \mathcal{X}_2^{(1)} \cup \mathcal{X}_2^{(2)} = \mathcal{X}$,然后对调二者来得到另一对数据集。我们再做第三次对折,因为每次对折我们得到两对数据,做 5 次对折我们共得到 10 个训练集和验证集:

$$\begin{aligned}\mathcal{T}_1 &= \mathcal{X}_1^{(1)} & \mathcal{V}_1 &= \mathcal{X}_1^{(2)} \\ \mathcal{T}_2 &= \mathcal{X}_1^{(2)} & \mathcal{V}_2 &= \mathcal{X}_1^{(1)} \\ \mathcal{T}_3 &= \mathcal{X}_2^{(1)} & \mathcal{V}_3 &= \mathcal{X}_2^{(2)} \\ \mathcal{T}_4 &= \mathcal{X}_2^{(2)} & \mathcal{V}_4 &= \mathcal{X}_2^{(1)} \\ &\vdots & & \\ \mathcal{T}_9 &= \mathcal{X}_5^{(1)} & \mathcal{V}_9 &= \mathcal{X}_5^{(2)} \\ \mathcal{T}_{10} &= \mathcal{X}_5^{(2)} & \mathcal{V}_{10} &= \mathcal{X}_5^{(1)}\end{aligned}$$

当然,我们可以做更多次对折以获得更多对的训练集和验证集,但是 Dietterich(1998)指出,在 5 次对折之后各集合共享了许多实例,过度的重叠使得由此计算的统计量,确切地说,验证误差率变得相互依赖而无法增加新的信息。即使 5 次对折,各集合也有重叠而统计

量也相互依赖,但是直到5折之前这些影响我们还是可以容忍的。而从另一方面来说,如果使用更少的对折次数,获得更少的数据(少于10组),将无法获得足够大的样本来拟合分布并进行假设检验。

19.6.3 自助法

为了从单个样本中产生多个样本,替代交叉验证的另一个选择是自助法(bootstrap),即从原始样本中以有放回地抽取实例的方法来产生新的样本。在17.6节我们看到,装袋使用自助法为不同的学习器产生训练集。自助样本可能比交叉验证样本有更多的重叠,因而其估计可能更为相互依赖;但对小数据集来说,这种方法被认为是做再抽样的最好方法。

在自助法中,我们从大小为 N 的数据集中有放回地抽取 N 个实例。原始数据集作为验证集。选取一个实例的概率为 $1/N$,不选取这个实例的概率为 $1 - 1/N$ 。一个实例在 N 次抽取均未选中的概率为:

$$\left(1 - \frac{1}{N}\right)^N \approx e^{-1} = 0.368$$

这意味着训练集包含了大约63.2%的实例;也就是说,系统未在36.8%的数据上进行训练,因而误差估计是悲观的。解决方法是重复该过程多次并观察平均行为。

19.7 度量分类器的性能

对于分类,特别是对于两类问题,已经提出了各种度量。存在4种情况,如表19-1所示。对于一个正实例,如果预测也是正的,则

它是一个真正(true positive)例;如果我们对正实例的预测是负的,则它是一个假负(false negative)例。对于一个负实例,如果预测也是负的,则它是一个真负(true negative)例;如果我们将负实例的预测为正的,则我们有一个假正(false positive)例。

在某些两类问题中,我们区分这两个类,因此有两种类型的错误:假正和假负。表19-2给出了适用于不同情况下的不同度量。让我们设想一种身份认证应用,其中用户通过声音登录他的账户。假正例是错误地允许冒名顶替者登录,而假负例是拒绝合法用户。显然,两种类型的错误并非同样糟糕,前一种更有害。真正率 $tp\text{-}rate$ 又称命中率(hit rate),度量通过身份认证的合法用户的比例,而假正率 $fp\text{-}rate$ 又称假报警率(false alarm rate),是错误地接受冒名顶替者的比例。

假设系统返回正类的概率 $\hat{P}(C_1 | x)$,而对于负类我们有 $\hat{P}(C_2 | x) = 1 - \hat{P}(C_1 | x)$,并且如果 $\hat{P}(C_1 | x) > \theta$,则我们选择“正”。如果 θ 接近于1,则我们很难选择正类;也就是说,

表 19-1 两类的混淆矩阵

实际的类	预测的类		
	正的	负的	合计
正的	tp : 真正	fn : 假负	p
负的	fp : 假正	tn : 真负	n
合计	p'	n'	N

表 19-2 两类问题使用的性能度量

名称	公式
误差 准确率	$(fp + fn) / N$ $(tp + tn) / N = 1 - \text{误差}$
$tp\text{-}rate$ $fp\text{-}rate$	tp / p fp / n
精度 召回率	tp / p' $tp / p = tp\text{-}rate$
灵敏度 特效性	$tp / p = tp\text{-}rate$ $tn / n = 1 - fp\text{-}rate$

我们将没有假正例，但是也只有少量真正例。随着减小 θ 来增加真正例数，我们冒着引进假正例的风险。

对于不同的 θ 值，我们可以得到许多 ($tp\text{-}rate$, $fp\text{-}rate$) 值对，而连接它们我们得到接受者操作特征(receiver operating characteristic, ROC)曲线，如图 19-3(a)所示。注意，不同的 θ 值对应于这两种类型误差的不同损失矩阵，并且 ROC 曲线也可以看作不同损失矩阵下的分类器的行为(参见习题 1)。

理想地，分类器的真正率为 1，假正率为 0，因此 ROC 曲线越靠近左上角的分类器越好。在对角线上，我们做出的真决策与假决策一样多，并且这是最坏情况(对角线下方的分类器都可以通过翻转它的决策来改进)。给定两个分类器，我们认为一个比另一个好，如果它的 ROC 曲线在另一个的上方；如果两条 ROC 曲线相交，则我们可以认为两个分类器在不同的损失条件下更好，如图 19-3(b)所示。

ROC 曲线提供了可视分析。如果我们想将该曲线归结为一个数，则可以通过计算曲线下方面积(area under the curve, AUC)来实现。理想地，分类器的 AUC 等于 1，并且可以比较不同分类器的 AUC 值，以得到不同损失条件下取平均的整体性能。

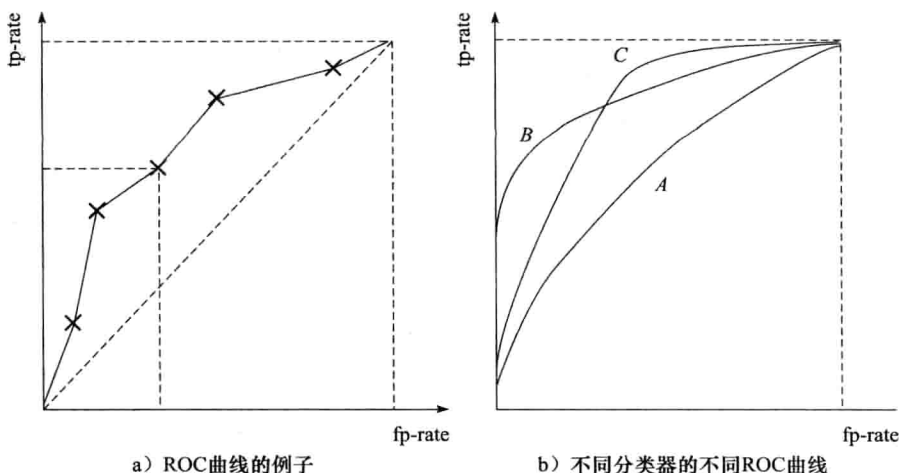


图 19-3 a) 典型的 ROC 曲线。每个分类器有一个阈值，使得我们可以在曲线上移动，并根据命中和假警报(即真正例和假正例)之间的相对重要程度来在曲线上确定一个点。ROC 曲线下方的面积称作 AUC。b) 如果一个分类器的 ROC 曲线更接近左上角(较大的 AUC)，则它更可取。B 和 C 都比 A 更可取；B 和 C 在不同的损失矩阵下更可取

在信息检索(information retrieval)中，有一个记录的数据库；例如，我们使用某些关键词提出查询，并且系统(基本上是一个两类分类器)返回大量记录。在该数据库中，存在一些相关记录，并且对于一个查询，系统可能检索到它们中的某些(真正例)，但可能不是所有的(假负例)；还可能错误地检索到不相关的记录(假正例)。相关和检索到的记录的集合可以用维恩图表示，如图 19-4(a)所示。精度(precision)是检索到的并且相关的记录数除以检索到的记录总数。如果精度为 1，则所有检索到的记录都是相关的，但可能还存在一些相关但未检索到的记录，如图 19-4(b)所示。召回率(recall)是检索到的并且相关的记录数除以相关记录的总数。即便召回率等于 1，所有相关记录可能都被检索到，但仍然可能有不相关的记录被检索到，如图 19-4(c)所示。与 ROC 曲线一样，我们也可以对不同的阈值绘制

精度和召回率的曲线。

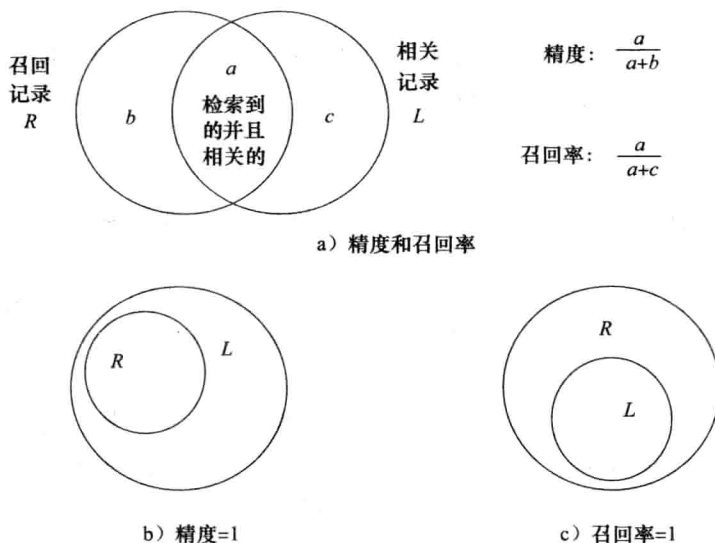


图 19-4 a) 使用维恩图定义精度和召回率。b) 精度为 1；所有检索到的记录都是相关的，但可能有相关的记录未被检索到。c) 召回率为 1；所有相关的记录都被检索到，但可能有不相关的记录也被检索到

492 从另一个角度但以相同的目的，还存在两个度量：灵敏度 (sensitivity) 和特效性 (specificity)。灵敏度与真正率和召回率相同。特效性度量检索负实例的好坏，它是真负实例数除以负实例的总数；它等于 1 减去假警报率。我们也可以使用不同的阈值绘制灵敏度和特效性的曲线。

对于 $K > 2$ 个类，如果我们使用 0/1 误差，则类混淆矩阵 (class confusion matrix) 是一个 $K \times K$ 的矩阵，其元素 (i, j) 是属于 C_i 类但却误分到 C_j 类的实例的个数。在理想情况下，所有的非对角线元素均应为 0，表示没有错误分类。类混淆矩阵允许我们准确地指出出现了哪种类型的误分类，即是否有两个类经常被混淆。或者，我们也可以分别定义 K 个两类问题，每个将一个类与其他 $K-1$ 个类分开。

19.8 区间估计

下面快速回顾一下我们将在假设检验中用到的区间估计 (interval estimation)。点估计，如最大似然估计，是对参数 θ 指定一个值。在区间估计中，我们以某种置信度对参数 θ 位于的区间进行确定。为了得到这种区间估计，我们利用点估计的概率分布。

例如，假设我们要从样本 $\mathbf{X} = \{x'_i\}_{i=1}^N$ 中估计正态密度的均值 μ 。 $m = \sum_i x'_i / N$ 是样本平均值，并且是对均值的点估计。 m 是正态分布值之和，因而也是正态的， $m \sim \mathcal{N}(\mu, \sigma^2/N)$ 。我们用单位正态分布来定义该统计量：

$$\frac{(m - \mu)}{\sigma / \sqrt{N}} \sim Z \quad (19.1)$$

我们知道 95% 的 Z 落在 $(-1.96, 1.96)$ 中，即 $P\{-1.96 < Z < 1.96\} = 0.95$ ，因而我们

有(参见图 19-5):

$$P\left\{-1.96 < \sqrt{N} \frac{(m - \mu)}{\sigma} < 1.96\right\} = 0.95$$

或等价地

$$P\left\{m - 1.96 \frac{\sigma}{\sqrt{N}} < \mu < m + 1.96 \frac{\sigma}{\sqrt{N}}\right\} = 0.95$$

493

也就是说“以 95% 的置信度”， μ 落在样本平均值正负 $1.96\sigma/\sqrt{N}$ 个单位的区间内，即双侧置信区间 (two-sided confidence interval)。以 99% 的置信度， μ 落在 $(m - 2.85\sigma/\sqrt{N}, m + 2.85\sigma/\sqrt{N})$ 中。也就是说，如果我们需要更高的置信度，则区间增大。随着样本集的规模 N 增大，区间变小。

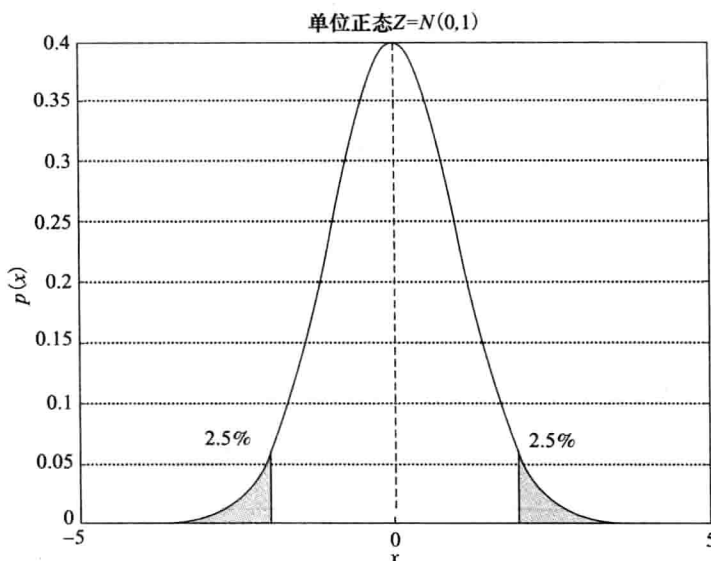


图 19-5 在单位正态分布中，95% 的值位于 $-1.96 \sim 1.96$

这可以按如下方法推广到任意置信度：令 z_α 使得

$$P\{Z > z_\alpha\} = \alpha, \quad 0 < \alpha < 1$$

由于 Z 是关于均值对称的，所以有 $z_{1-\alpha/2} = -z_{\alpha/2}$ ，并且 $P\{X < -z_{\alpha/2}\} = P\{X > z_{\alpha/2}\} = \alpha/2$ 。因而，对于任意给定的置信水平 $1 - \alpha$ ，我们有：

$$P\{-z_{\alpha/2} < Z < z_{\alpha/2}\} = 1 - \alpha$$

并且

$$P\left\{-z_{\alpha/2} < \sqrt{N} \frac{(m - \mu)}{\sigma} < z_{\alpha/2}\right\} = 1 - \alpha$$

494

或

$$P\left\{m - z_{\alpha/2} \frac{\sigma}{\sqrt{N}} < \mu < m + z_{\alpha/2} \frac{\sigma}{\sqrt{N}}\right\} = 1 - \alpha \quad (19.2)$$

因而，对于任意 α ，可以对 μ 计算置信度为 $100(1 - \alpha)\%$ 的双侧置信区间。

类似地，如有 $P\{Z < 1.64\} = 0.95$ ，则我们有(参见图 19-6)

$$P\left\{\sqrt{N}\frac{(m-\mu)}{\sigma} < 1.64\right\} = 0.95$$

或

$$P\left\{m - 1.64\frac{\sigma}{\sqrt{N}} < \mu\right\} = 0.95$$

并且 $(m - 1.64\sigma/\sqrt{N}, \infty)$ 是 μ 的 95% 的单侧上置信区间 (one-sided upper confidence interval), 其定义了一个下界。一般而言, μ 的 $100(1-\alpha)\%$ 的单侧置信区间可通过下式计算:

$$P\left\{m - z_{\alpha}\frac{\sigma}{\sqrt{N}} < \mu\right\} = 1 - \alpha \quad (19.3)$$

类似地, 可以计算定义上界的单侧下置信区间。

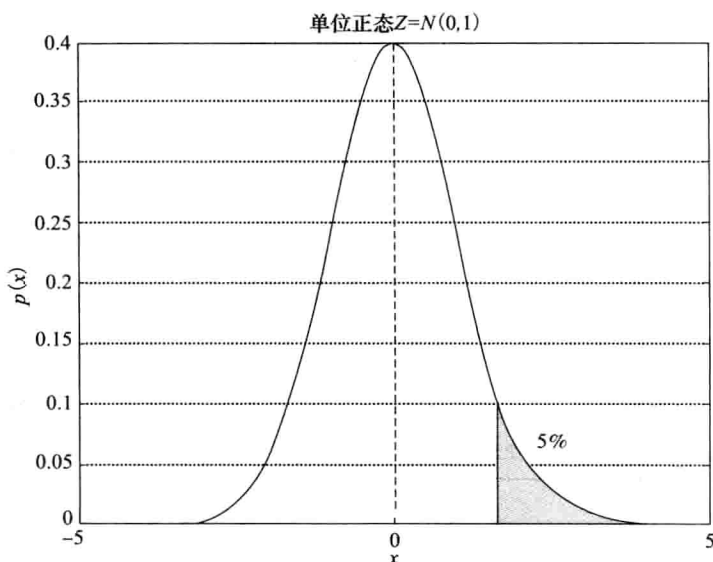


图 19-6 单位正态分布的 95% 落在 1.64 之前

在上述各区间中, 我们使用了 σ , 即我们假定方差是已知的。如果方差未知, 则我们可用样本方差

$$S^2 = \sum (x' - m)^2 / (N - 1)$$

来替代 σ^2 。我们知道当 $x' \sim \mathcal{N}(\mu, \sigma^2)$ 时, $(N-1)S^2/\sigma^2$ 是自由度为 $N-1$ 的卡方(分布)。我们同样知道 m 和 S^2 是相互独立的。于是, $\sqrt{N}(m-\mu)/S$ 是自由度为 $N-1$ 的 t 分布(参见 A. 3.7 节), 记作

$$\frac{\sqrt{N}(m-\mu)}{S} \sim t_{N-1} \quad (19.4)$$

因而, 对任意的 $\alpha \in (0, 1/2)$, 我们可以使用该 t 分布 (t distribution) 而非单位正态分布 Z 确定的值来定义一个区间:

$$P\left\{t_{1-\alpha/2, N-1} < \sqrt{N}\frac{(m-\mu)}{S} < t_{\alpha/2, N-1}\right\} = 1 - \alpha$$

或使用 $t_{1-\alpha/2, N-1} = -t_{\alpha/2, N-1}$, 我们有

$$P\left\{m - t_{\alpha/2, N-1} \frac{S}{\sqrt{N}} < \mu < m + t_{\alpha/2, N-1} \frac{S}{\sqrt{N}}\right\} = 1 - \alpha$$

类似地, 可以定义单侧置信区间。 t 分布比单位正态分布有更大的展宽(较长的尾), 因而 t 分布给出的区间一般更大, 考虑到未知方差所引入的附加的不确定性的存在, 这应该在预料之中。

19.9 假设检验

在某些应用中, 我们可能希望使用样本对涉及参数的一些特定假设进行检验, 而不是显式地估计某些参数。例如, 我们可能希望检验均值是否小于 0.02, 而不是估计均值。如果随机样本与所考虑的假设一致, 则我们“接受”该假设, 否则, 我们说它“被拒绝”。但是, 当我们做这样的决定时, 实际上并非说假设为真或假, 而是说在一定的置信程度上, 样本数据和假设看起来是一致的。

在假设检验(hypothesis testing)中, 方法如下: 我们定义一个统计量, 如果假设正确, 则该统计量服从某一分布。如果从样本中计算的统计量具有很低的概率是取自该分布, 则我们拒绝该假设; 否则, 我们无法拒绝它。

假设有一个样本取自一个均值 μ 未知、方差 σ^2 已知的正态分布, 而我们对一个关于 μ 的假设进行检验, 例如, 其值是否等于一个指定的常数 μ_0 。该假设记作 H_0 并称为原假设(null hypothesis)

$$H_0: \mu = \mu_0$$

相对的备择假设为:

$$H_1: \mu \neq \mu_0$$

m 是 μ 的点估计, 而且如果 m 距离 μ_0 太远, 则拒绝 H_0 是合理的, 这正是要使用区间估计的地方。我们以显著水平(level of significance) α 无法拒绝该假设, 如果 μ_0 位于 $100(1 - \alpha)\%$ 的置信区间, 即如果

$$\frac{\sqrt{N}(m - \mu_0)}{\sigma} \in (-z_{\alpha/2}, z_{\alpha/2}) \quad (19.5)$$

如果它落在任何一侧的外面, 则我们拒绝原假设。这是一个双侧检验(two-sided test)。

如果假设是正确时我们拒绝了它, 这就是一个第一类错误(type I error), 而在检验之前设定的 α 值定义了我们可以多大程度上容忍第一类错误, α 的通常取值为 0.1、0.05、0.01(参见表 19-3)。第二类错误(type II error)是如果真实均值 μ 不等于 μ_0 时我们无法拒绝原假设。当真实均值为 μ 时不拒绝 H_0 的概率是 μ 的函数, 并由下式给出:

$$\beta(\mu) = P_{\mu}\left\{-z_{\alpha/2} \leq \frac{m - \mu_0}{\sigma/\sqrt{N}} \leq z_{\alpha/2}\right\} \quad (19.6)$$

$1 - \beta(\mu)$ 称为检验的功效函数(power function), 并且等于当 μ 为真实值时假设被拒绝的概率。随着 μ 与 μ_0 接近, 第二类错误的概率增加, 并且我们可以计算为了能够以足够的功

表 19-3 第一类错误、第二类错误和检验功效

事实	决策	
	无法拒绝	拒绝
真	正确	第一类错误
假	第二类错误	正确(功效)

效检测出差 $\delta = |\mu - \mu_0|$ ，需要多大的样本。

与备择假设为 $\mu \neq \mu_0$ 时的双侧检验相反，也可以进行如下形式的单侧检验 (one-sided test)：

$$H_0: \mu \leq \mu_0, H_1: \mu > \mu_0$$

显著性水平为 α 的单侧检验定义了界定于单侧的 $100(1 - \alpha)\%$ 置信区间，为了接受假设， m 的值必须落在该区间内。我们无法拒绝该假设，如果

$$\frac{\sqrt{N}}{\sigma}(m - \mu_0) \in (-\infty, z_\alpha) \quad (19.7)$$

并且落在外边时拒绝。注意，原假设 H_0 也允许等式，这意味着当检验被拒绝时我们只得到序信息。这告诉我们应该使用哪个单侧检验。我们的任何断言都应该在 H_1 中，使得检验的拒绝将支持我们的断言。

如果方差未知，我们可以像在区间估计中所作的那样，以样本方差来替代总体方差并利用下述事实：

$$\frac{\sqrt{N}(m - \mu_0)}{S} \sim t_{N-1} \quad (19.8)$$

例如，对于 $H_0: \mu = \mu_0$ 和 $H_1: \mu \neq \mu_0$ ，我们以显著性水平 α 接受假设，如果

$$\frac{\sqrt{N}(m - \mu_0)}{S} \in (-t_{\alpha/2, N-1}, t_{\alpha/2, N-1}) \quad (19.9)$$

这就是双侧 t 检验 (two-sided t test)。单侧 t 检验可以类似地定义。

19.10 评估分类算法的性能

既然我们已经回顾了假设检验，我们看看如何将其应用于错误率检验。我们将讨论分类误差，但是只要我们能够为抽样分布确定适当的参数形式，同样的技术就可以用于回归的均方误差、非监督学习的对数似然、增强学习的期望奖励等，我们还将讨论当找不到这种参数形式时的非参数检验。

我们从错误率评估入手，在下一节讨论错误率比较。

19.10.1 二项检验

我们从只有一个训练集 \mathcal{T} 和一个验证集 \mathcal{V} 的情形入手。我们在 \mathcal{T} 上训练分类器并在 \mathcal{V} 上检验它。我们以 p 表示分类器产生一个误分类错误的概率。我们不知道 p ，要对它进行估计或对关于它的假设进行检验。对于来自验证集 \mathcal{V} 的索引号为 t 的实例，令 x^t 表示分类器决策的正确性： x^t 是一个 0/1 伯努利随机变量，当分类器产生一次错误时它取值 1，而当分类器正确时它取值 0。二项随机变量 X 表示错误的总数：

$$X = \sum_{t=1}^N x^t$$

我们想检验错误的概率 p 是否小于或等于我们指定的某个值 p_0 ：

$$H_0: p \leq p_0 \text{ 与 } H_1: p > p_0$$

如果错误的概率为 p ，则分类器在 N 次分类中犯 j 次错误的概率为：

$$P\{X = j\} = C_N^j p^j (1 - p)^{N-j}$$

如果我们看到 $X = e$ 个或更多错误的概率很小, 则拒绝 $p \leq p_0$ 是合理的。即二项检验 (binomial test) 拒绝该假设, 如果

$$P\{X \geq e\} = \sum_{x=e}^N C_N^x p_0^x (1 - p_0)^{N-x} < \alpha \quad (19.10)$$

其中, α 是显著性, 例如为 0.05。

19.10.2 近似正态检验

如果 p 是错误的概率, 则我们的点估计是 $\hat{p} = X/N$ 。于是, 如果 \hat{p} 比 p_0 大很多, 则拒绝原假设是合理的。多大才算足够大由 \hat{p} 的选择分布和显著性 α 给定。

因为 X 是服从相同分布的独立随机变量之和, 依据中心极限定理, 对于大的 N 值, X/N 近似服从均值为 p_0 、方差为 $p_0(1 - p_0)$ 的正态分布。于是

$$\frac{X/N - p_0}{\sqrt{p_0(1 - p_0)/N}} \sim Z \quad (19.11)$$

其中 \sim 表示“近似服从分布”。于是, 使用式 (19.7), 当 $X = e$ 时, 如果上式的值大于 z_α , 则近似正态检验 (approximate normal test) 拒绝原假设。 $z_{0.05}$ 的值是 1.64。只要 N 的值不是太小, 并且 p 的值不是非常接近于 0 或 1, 这个近似将很有效; 作为一种经验法则, 我们要求 $Np \geq 5$, 并且 $N(1 - p) \geq 5$ 。

19.10.3 t 检验

前面讨论的两种检验方法都使用一个验证集。如果在 K 对训练集和验证集上运行算法 K 次, 则我们在 K 个验证集上得到 K 个错误百分比 p_i , $i = 1, \dots, K$ 。如果在 \mathcal{T}_i 上训练的分类器对 \mathcal{V}_i 中的实例 t 产生了一次误分, 则令 x_i^t 为 1; 否则令 x_i^t 为 0。于是

$$p_i = \frac{\sum_{t=1}^N x_i^t}{N}$$

由于有

$$m = \frac{\sum_{i=1}^K p_i}{K}, \quad S^2 = \frac{\sum_{i=1}^K (p_i - m)^2}{K - 1}$$

根据式 (19.8), 我们有

$$\frac{\sqrt{K}(m - p_0)}{S} \sim t_{K-1} \quad (19.12)$$

如果上式的值大于 $t_{\alpha, K-1}$, 则 t 检验拒绝“分类算法以显著性水平 α 具有 p_0 或更低的错误率”的原假设。通常, K 取值为 10 或 30。 $t_{0.05, 9} = 1.83$ 而 $t_{0.05, 29} = 1.70$ 。

19.11 比较两个分类算法

给定两个学习算法和一个训练集, 我们想要比较和检验这两个算法所构建的分类器是否

具有相同的期望错误率。

19.11.1 McNemar 检验

给定一个训练集和一个验证集，我们使用两个算法在训练集上训练两个分类器，在验证集上对它们进行检验并计算它们的误差。如下所示的相依表(contingency table)是一个矩阵形式的自然数数组，表示各种情况的计数或频率：

e_{00} ：两个分类器都错误分类的实例个数	e_{01} ：分类器 1 错误分类而分类器 2 没有错误分类的实例个数
e_{10} ：分类器 2 错误分类而分类器 1 没有错误分类的实例个数	e_{11} ：两个分类器都正确分类的实例个数

在两个分类算法有相同的错误率的原假设下，我们期望 $e_{01} = e_{10}$ ，并且二者都等于 $(e_{01} + e_{10})/2$ 。我们有自由度为 1 的卡方统计量

$$\frac{(|e_{01} - e_{10}| - 1)^2}{e_{01} + e_{10}} \sim \chi_1^2 \quad (19.13)$$

并且如果这个值大于 $\chi_{\alpha,1}^2$ ，则 McNemar 检验拒绝两个分类算法以显著水平 α 具有相同错误率的假设。对于 $\alpha = 0.05$ ， $\chi_{0.05,1}^2 = 3.84$ 。

19.11.2 K-折交叉验证配对 t 检验

使用 K-折交叉验证在数据集上产生 K 对训练集和验证集。我们使用两个分类算法在训练集 $\mathcal{T}_i (i=1, \dots, K)$ 上训练并在验证集 \mathcal{V}_i 上检验。两个分类器在验证集上的误差率分别记作 p_i^1 和 p_i^2 。

如果两个分类算法具有相同的错误率，则我们预期它们具有相同的均值，或等价地说，它们的均值之差为 0。在第 i 折，两个分类器的错误率之差是 $p_i = p_i^1 - p_i^2$ 。这是配对检验(paired test)，即对于每个 i，两个算法都使用相同的训练集和验证集。K 次比较之后，得到一个包含 K 个点的 p_i 分布。假定 p_i^1 和 p_i^2 都是(近似)正态的，则其差 p_i 也是正态的。原假设为该分布的均值为 0：

$$H_0: \mu = 0, \quad H_1: \mu \neq 0$$

我们定义

$$m = \frac{\sum_{i=1}^K p_i}{K}, \quad S^2 = \frac{\sum_{i=1}^K (p_i - m)^2}{K-1}$$

在 $\mu=0$ 的原假设下，我们有一个统计量，它是自由度为 $K-1$ 的 t 分布：

$$\frac{\sqrt{K}(m-0)}{S} = \frac{\sqrt{K} \cdot m}{S} \sim t_{K-1} \quad (19.14)$$

因而，如果该值落在区间 $(-t_{\alpha/2, K-1}, t_{\alpha/2, K-1})$ 之外，则 K-折交叉验证配对 t 检验(K-fold cv paired t test)拒绝两个分类算法以显著水平 α 具有相同误差率的假设。 $t_{0.025,9} = 2.26$ ， $t_{0.025,29} = 2.05$ 。

如果我们想检验第一个算法的错误率是否比第二个算法更小,则需要使用单侧假设,并使用单尾检验:

$$H_0: \mu \geq 0 \text{ 与 } H_1: \mu < 0$$

如果检验拒绝,则我们断言第一个算法具有显著较小的错误率就得到了支持。

19.11.3 5×2 交叉验证配对 t 检验

在 Dietterich(1998)提出的 5×2 交叉验证 t 检验中,我们进行五轮对折交叉验证。在每一轮中,数据集被划分为两个大小相等的集合。 $p_i^{(j)}$ 表示两个分类器在第 i 轮中第 j 折数据上的错误率之差,其中 $i=1, \dots, 5, j=1, 2$ 。在第 i 轮的平均为 $\bar{p}_i = (p_i^{(1)} + p_i^{(2)})/2$, 估计方差为 $s_i^2 = (p_i^{(1)} - \bar{p}_i)^2 + (p_i^{(2)} - \bar{p}_i)^2$ 。

在两个分类算法具有相同错误率的原假设下, $p_i^{(j)}$ 是两个同分布的比例值之差,而忽略这些比例值不相互独立的事实, $p_i^{(j)}$ 可近似地被认为服从均值为 0, 方差 σ^2 未知的正态分布。于是, $p_i^{(j)}/\sigma$ 是近似单位正态的。如果假定 $p_i^{(1)}$ 和 $p_i^{(2)}$ 是独立、正态的(严格地说并非如此,因为它们的训练和验证集并非相互独立抽取的), 则 s_i^2/σ^2 服从自由度为 1 的卡方分布。如果假定每个 s_i^2 是相互独立的(事实并非如此,因为它们从相同的可用数据集计算得到), 则它们的和服从自由度为 5 的卡方分布:

$$M = \frac{\sum_{i=1}^5 s_i^2}{\sigma^2} \sim \chi_5^2$$

并且

$$t = \frac{p_1^{(1)}}{\sqrt{M/5}} = \frac{p_1^{(1)}}{\sqrt{\sum_{i=1}^5 s_i^2/5}} \sim t_5 \quad (19.15)$$

给出了自由度为 5 的 t 统计量。如果该统计量的值落在区间 $(-t_{\alpha/2,5}, t_{\alpha/2,5})$ 之外, 则 5×2 交叉验证 t 检验 (5×2 cv paired t test) 拒绝两个分类算法以显著性水平 α 具有相同的错误率的假设。 $t_{0.025,5} = 2.57$ 。

19.11.4 5×2 交叉验证配对 F 检验

我们注意到式(19.15)中的分子 $p_1^{(1)}$ 是任意的; 实际上, 有 10 个不同的值可当作分子, 即 $p_i^{(j)}, j=1, 2, i=1, \dots, 5$, 产生 10 个可能的统计量:

$$t_i^{(j)} = \frac{p_i^{(j)}}{\sqrt{\sum_{i=1}^5 s_i^2/5}} \quad (19.16)$$

Alpaydin(1999)提出了 5×2 交叉验证 t 检验的扩展, 组合 10 个可能的统计量的结果。如果 $p_i^{(j)}/\sigma \sim Z$, 则有 $(p_i^{(j)})^2/\sigma^2 \sim \chi_1^2$, 并且它们的和是自由度为 10 的卡方分布:

$$N = \frac{\sum_{i=1}^5 \sum_{j=1}^2 (p_i^{(j)})^2}{\sigma^2} \sim \chi_{10}^2$$

将上式作为式(19.15)的分子,得到的统计量为两个卡方分布随机变量的比值。两个变量分别除以它们的自由度得到第一自由度为 10,第二自由度为 5 的 F 分布(参见 A.3.8 节):

$$f = \frac{N/10}{M/5} = \frac{\sum_{i=1}^5 \sum_{j=1}^2 (p_i^{(j)})^2}{2 \sum_{i=1}^5 s_i^2} \sim F_{10,5} \quad (19.17)$$

如果该值大于 $F_{\alpha,10,5}$, 则 5×2 交叉验证配对 F 检验 (5×2 cv paired F test) 拒绝两个分类算法以显著性水平 α 具有相同错误率的假设。 $F_{0.05,10,5} = 4.74$ 。

19.12 比较多个算法: 方差分析

在很多情况下,有多个算法,我们希望比较它们的期望误差。给定 L 个算法,我们在 K 个训练集上对它们进行训练,每个算法产生 K 个分类器,而后在 K 个验证集上进行检验并记录相应的错误率。这样产生了 L 组,每组 K 个值。于是,问题是比较这 L 个样本差异的统计显著性。这是一个具有单个因素、 L 个水平(学习算法)的实验,对每个水平重复 K 次。

在方差分析 (analysis of variance, ANOVA) 中,我们考虑 L 个独立的样本,每个大小为 K ,由未知均值 μ_j 和未知公共方差 σ^2 的正态随机变量组成:

$$X_{ij} \sim \mathcal{N}(\mu_j, \sigma^2), j = 1, \dots, L, i = 1, \dots, K$$

我们想对“所有均值相等”的假设 H_0 进行检验:

$$H_0: \mu_1 = \mu_2 = \dots = \mu_L \text{ 与 } H_1: \text{至少在一对}(r, s) \text{ 上, } \mu_r \neq \mu_s$$

对多个分类算法的错误率进行比较就属于这种情况。我们有 L 个分类算法,并且有其在 K 个验证折上的错误率。 X_{ij} 是分类算法 j 在第 i 折数据所训练的分类器的验证错误数。每个 X_{ij} 都是二项的并且是近似正态的。如果 H_0 未被拒绝,则我们在这 L 个分类算法的误差率之间找不到显著差别。因而,这是我们在 19.11 节所看到的对两个分类算法进行比较的检验方法的推广。这 L 个分类算法可能不同或使用不同的超参数,如多层感知器的隐藏单元数, k -nn 的近邻数等。

方差分析方法导出两个关于 σ^2 的估计。第一个估计只有在 H_0 为真的时候才为真,而第二个估计始终是一个有效估计,无论 H_0 是否为真。如果两个估计显著不同,则方差分析拒绝 H_0 ,即拒绝 L 个样本取自相同的总体。

第一个关于 σ^2 的估计有效,仅当假设为真,即仅当 $\mu_j = \mu, j = 1, \dots, L$ 。如果 $X_{ij} \sim \mathcal{N}(\mu, \sigma^2)$, 则组平均

$$m_j = \sum_{i=1}^K \frac{X_{ij}}{K}$$

也是正态的,均值为 μ 、方差为 σ^2/K 。如果假设为真,则 $m_j (j = 1, \dots, L)$ 是 L 个取自 $\mathcal{N}(\mu, \sigma^2/K)$ 的实例。于是,它们的均值和方差分别为:

$$m = \frac{\sum_{j=1}^L m_j}{L}, \quad S^2 = \frac{\sum_{j=1}^L (m_j - m)^2}{L-1}$$

因而, σ^2 的一个估计是 $K \cdot S^2$, 即

$$\hat{\sigma}_b^2 = K \sum_{j=1}^L \frac{(m_j - m)^2}{L - 1} \quad (19.18)$$

每个 m_j 都是正态的，并且 $(L-1)S^2/(\sigma^2/K)$ 是自由度为 $(L-1)$ 的卡方分布。因而，我们有

$$\sum_j \frac{(m_j - m)^2}{\sigma^2/k} \sim \chi_{L-1}^2 \quad (19.19)$$

我们定义组间平方和 SS_b 为

$$SS_b \equiv K \sum_j (m_j - m)^2$$

因而，当 H_0 为真时，我们有

$$\frac{SS_b}{\sigma^2} \sim \chi_{L-1}^2 \quad (19.20)$$

第二个关于 σ^2 的估计是组方差 S_j^2 的平均值。组方差定义为

$$S_j^2 = \frac{\sum_{i=1}^K (X_{ij} - m_j)^2}{K - 1}$$

而它们的平均值为：

$$\hat{\sigma}^2 = \sum_{j=1}^L \frac{S_j^2}{L} = \sum_j \sum_i \frac{(X_{ij} - m_j)^2}{L(K-1)} \quad (19.21)$$

我们定义组内平方和 SS_w 为

$$SS_w \equiv \sum_j \sum_i (X_{ij} - m_j)^2$$

回忆一下，对正态样本，我们有

$$(K-1) \frac{S_j^2}{\sigma^2} \sim \chi_{K-1}^2$$

并且卡方分布之和仍然是卡方分布，我们有

$$(K-1) \sum_{j=1}^L \frac{S_j^2}{\sigma^2} \sim \chi_{L(K-1)}^2$$

因而

$$\frac{SS_w}{\sigma^2} \sim \chi_{L(K-1)}^2 \quad (19.22)$$

于是，我们的任务是比较两个方差是否相等，这可以通过检查它们的比值是否接近于 1 来实现。两个独立卡方随机变量分别除以其相应的自由度的比值是一个服从 F 分布的随机变量，因而当 H_0 为真时，我们有

$$F_0 = \left(\frac{SS_b/\sigma^2}{L-1} \right) / \left(\frac{SS_w/\sigma^2}{L(K-1)} \right) = \frac{SS_b/(L-1)}{SS_w/L(K-1)} = \frac{\hat{\sigma}_b^2}{\hat{\sigma}_w^2} \sim F_{L-1, L(K-1)} \quad (19.23)$$

对于任意给定的显著性水平值 α ，如果该统计量的值大于 $F_{\alpha, L-1, L(K-1)}$ ，则拒绝 L 个分类算法具有相同的期望错误率的假设。

注意，如果两个估计显著不一致，则我们拒绝 H_0 。如果 H_0 不成立，则 m_j 在 m 附近的

方差一般比 H_0 成立时大。因此, 如果 H_0 不成立, 则第一个估计 $\hat{\sigma}_b^2$ 将是 σ^2 的过高估计, 而该比值将大于 1。对于 $\alpha = 0.05$, $L = 5$ 和 $K = 10$, $F_{0.05, 4, 45} = 2.6$ 。如果 X_{ij} 在 m 附近以方差 σ^2 变化, 那么如果 H_0 成立, 则 m_j 在 m 附近将以方差 σ^2/K 变化。看来, 如果它们变化很大, 则应该拒绝 H_0 , 因为 m_j 在 m 附近的位移很难再用某种不断增加的噪声解释。

方差分析的名字源于将数据中的总变化划分成它的分量。

$$SS_T \equiv \sum_j \sum_i (X_{ij} - m)^2 \quad (19.24)$$

SS_T 除以它的自由度 $K \cdot L - 1$ (存在 $K \cdot L$ 个数据点, 我们损失了一个自由度, 因为 m 是固定的) 给出 X_{ij} 的样本方差。可以证明 (习题 5), 总平方和可以分解成组间平方和与组内平方和

$$SS_T = SS_b + SS_w \quad (19.25)$$

ANOVA (方差分析) 的结果在如表 19-4 所示的 ANOVA 表中。这是方差的基本一路 (one-way) 分析, 其中只有一个因素, 例如学习算法。我们可以考虑用多个因素实验, 例如, 一个因素可以是分类算法, 另一个因素是分类之前使用的特征提取算法, 而这是一个具有相互影响的双因素实验 (two-factor experiment with interaction)。

如果假设被拒绝, 我们只知道 L 个分组之间存在某种差异, 但是我们并不知道差异在何处。为此, 我们做事后检验 (posthoc testing), 即一组涉及分组的子集的额外检验, 例如逐对。

表 19-4 单个因素模型的方差分析 (ANOVA) 表

变化源	平方和	自由度	均方	F_0
组间	$SS_b \equiv K \sum_j (m_j - m)^2$	$L - 1$	$MS_b = \frac{SS_b}{L - 1}$	$\frac{MS_b}{MS_w}$
组内	$SS_w \equiv \sum_j \sum_i (X_{ij} - m_j)^2$	$K(L - 1)$	$MS_w = \frac{SS_w}{L(K - 1)}$	
总和	$SS_T \equiv \sum_j \sum_i (X_{ij} - m)^2$	$L \cdot K - 1$		

费希尔的最小方差检验 (least square difference test, LSD) 以逐对的方式比较分组。对于每个分组, 我们有 $m_i \sim \mathcal{N}(\mu_i, \sigma_w^2 = MS_w/K)$, $m_i - m_j \sim \mathcal{N}(\mu_i - \mu_j, 2\sigma_w^2)$ 。于是, 在原假设 $H_0: \mu_i = \mu_j$ 下, 我们有

$$t = \frac{m_i - m_j}{\sqrt{2\sigma_w^2}} \sim t_{L(K-1)}$$

如果 $|t| > t_{\alpha/2, L(K-1)}$, 则我们拒绝 H_0 , 支持备择假设 $H_1: \mu_i \neq \mu_j$ 。类似地, 可以定义单边检验, 以找出逐对的顺序。

当我们做大量实验以提取结论时, 这称为多重比较 (multiple comparisons), 并且需要记住: 如果我们以显著性水平 α 对 T 个假设进行检验, 则至少有一个假设不正确地被拒绝的概率至多为 $T\alpha$ 。例如, 均以 95% 的单个置信区间进行计算的 6 个置信区间同时正确的概率至少为 70%。因而, 为了确保整体置信区间至少为 $100(1 - \alpha)$, 单个置信区间应当设置为 $100(1 - \alpha/T)$, 这称为 Bonferroni 校正 (Bonferroni correction)。

有时, 可能出现这种情况: ANOVA 拒绝了原假设, 而事后逐对检验都找不到显著的差

别。在这种情况下，我们的结论是：均值之间存在差别，但是我们需要更多的数据，以便指出差别到底在何处。

注意，分析的主要开销是在 K 个训练集和验证集上对 L 个分类算法进行训练和检验。一旦训练、检验完成并且结果存储在一张 $K \times L$ 的表中，则从表中进行方差分析或计算两两比较检验统计量的开销就相对很小。

19.13 在多个数据集上比较

假定我们想在多个数据集而不是在一个数据集上比较两个或多个算法。这与前面的不同之处在于：依赖于其归纳偏倚与问题匹配程度的算法在不同的数据集上将会不同，并且不能再认为不同数据集上的误差值在平均准确率附近是正态分布的。这意味着我们上一节讨论的基于二项分布是近似正态分布的参数检验不再能用，而我们需要借助于非参数检验 (non-parametric test)。这种检验的优点是，我们还可以使用它们比较其他非正态的统计量，例如，训练时间、自由参数的个数等。

对于稍微偏离正态分布的情况，参数检验一般是鲁棒的，特别是当样本很大时。非参数检验不受分布的限制，但不太有效；即如果两种检验都可以使用，则优先选择参数检验。为了取得相同的功效，对应的非参数检验需要更大的样本。非参数检验并不假定知道基本总体的分布，而只假定值可以比较或是有序的，并且正如我们将看到，非参数检验使用这种序信息。

当我们将一个算法在许多不同的数据集上训练后，该算法在这些数据集上误差的平均值不是一个有意义的值，并且我们不能使用这种值，例如，比较两个算法 A 和 B 。为了比较两个算法，唯一可以使用的信息是，是否在任意数据集上 A 比 B 准确；于是，我们可以统计 A 比 B 准确的次数，并检查这是否可能是偶然的，它们是否确实同样准确。当有更多算法时，我们将考察被不同算法训练的学习器的平均等级 (rank)。非参数检验基本上使用这些等级数据，而非绝对的数值。

在继续深入讨论这些检验之前，应该强调的是，在各种各样的应用上比较算法的错误率是没有意义的。因为不存在“最佳学习算法”这种事，因此这种检验不会是令人信服的。然而，我们可以在相同应用的大量数据集或版本上比较算法。例如，对于人脸识别，我们可以有许多不同的数据集，但具有不同的性质 (分辨率、光照、研究对象数等)，而我们可以使用非参数检验在这些数据集上比较算法。数据集的不同性质使得我们不可能将来自不同数据集的图像混在一起，形成单个数据集，但是我们可以在不同的数据集上分别训练算法，分别得到等级，并组合这些等级得到总体决策。

19.13.1 比较两个算法

假定我们想比较两个算法。我们采用逐对方式在 N 个不同的数据集上训练和验证它们；即除不同的算法之外，所有的条件都应该是一样的。我们得到结果 e_i^1 和 e_i^2 ，并且如果我们在每个数据集上都使用 K 折交叉验证，则这些是 K 个值的均值或中位数。符号检验 (sign test) 基于如下思想：如果两个算法具有相同的误差，则在每个数据集上，第一个算法的误差比第二小的概率为 $1/2$ ，因而我们预料第一个算法在 $N/2$ 个数据集上获胜，我们定义

$$X_i = \begin{cases} 1 & \text{如果 } e_i^1 < e_i^2 \text{ 并且 } X = \sum_{i=1}^N X_i \\ 0 & \text{否则} \end{cases}$$

假定我们想检验

$$H_0: \mu_1 \geq \mu_2 \text{ 与 } H_1: \mu_1 < \mu_2$$

如果原假设是正确的, 则 X 在 N 次试验上是二项的, 其中 $p = 1/2$ 。假设第一个算法在 $X = e$ 个数据集上获胜。于是, 当确实有 $p = 1/2$ 时, 获胜的次数为 e 或更少的概率为

$$P\{X \leq e\} = \sum_{x=0}^e C_N^x \left(\frac{1}{2}\right)^x \left(\frac{1}{2}\right)^{N-x}$$

如果这个概率值太小, 即小于 α , 则我们拒绝原假设。如果出现平局, 则我们将它们同等地划归两边; 也就是说, 如果有 t 次平局, 则我们将 $t/2$ 加到 e 上(如果 t 是奇数, 则我们忽略一次, 并将 N 减 1)。

在检验

$$H_0: \mu_1 \leq \mu_2 \text{ 与 } H_1: \mu_1 > \mu_2$$

中, 如果 $P\{X \geq e\} < \alpha$, 则我们拒绝 H_0 。

对于双边检验

$$H_0: \mu_1 = \mu_2 \text{ 与 } H_1: \mu_1 \neq \mu_2$$

如果 e 太小或太大, 则我们拒绝原假设。当 $e < N/2$ 时, 如果 $2P\{X \leq e\} < \alpha$, 则我们拒绝原假设; 当 $e > N/2$ 时, 如果 $2P\{X \geq e\} < \alpha$, 则我们拒绝原假设——我们需要找出对应的尾并将其乘以 2, 因为这是双边检验。

正如我们在前面所讨论的, 非参数检验可以用来比较任意度量, 例如训练时间。在这种情况下, 我们看到非参数学习的优点: 它使用序, 而不是绝对的平均值。假设我们在 10 个数据集上比较两个算法, 其中 9 个数据集很小, 并且两个算法的训练时间都以分钟计, 而一个数据集很大, 其训练时间以天计。如果我们使用参数检验, 并取训练时间的平均值, 则一个大数据集就左右了决策; 而当我们使用非参数检验并分别在每个数据集上比较时, 使用序具有在每个数据集上规范化的效果, 因此有助于我们得到鲁棒的决策。

我们也可将符号检验用作一个样本的检验, 例如, 不是将 μ_1 与第二个总体的均值比较, 而是将它与一个常数 μ_0 比较, 来检查算法在所有数据集上的平均误差是否小于 2%。我们可以简单地通过如下方法来做这件事: 将常量 μ_0 取代第二个样本的所有观测, 并使用前面的过程; 即我们将统计误差大于和小于 0.02 的次数, 并检查在原假设下这是否太不可能。对于较大的 N , 可以使用正态分布来近似二项分布(习题 6), 但实践中, 数据集的个数可能小于 20。注意, 符号检验是在总体中位数上的检验, 如果分布是对称的, 则中位数等于均值。

符号检验只使用差的符号, 而不是差的量值, 但是我们可能面对这样一种情况: 当第一个算法获胜时, 总是大幅度赢; 而当第二个算法获胜时, 总是勉强赢。Wilcoxon 符号秩检验 (Wilcoxon signed rank test) 同时使用符号和差的量值, 其方法如下:

假设除了差的符号之外, 我们还计算 $m_i = |e_i^1 - e_i^2|$, 并且将它们排序, 使得最小的 m_i 的秩为 1, 次最小的秩为 2, 以此类推。如果出现平局, 则它们的秩取它们稍微不同时得到的平均值。例如, 如果这些差值为 2、1、2、4, 则秩为 2.5、1、2.5、4。然后, 我们计算

符号为正的所有秩的和 w_+ , 符号为负的所有秩的和 w_- 。

仅当 w_+ 远小于 w_- 时, 可以拒绝原假设 $\mu_1 \leq \mu_2$, 而取备择假设 $\mu_1 > \mu_2$ 。类似地, 仅当 w_+ 和 w_- 都很小, 即当 $w = \min(w_+, w_-)$ 很小时, 可以拒绝双边假设 $\mu_1 = \mu_2$, 而接受备择假设 $\mu_1 \neq \mu_2$ 。Wilcoxon 符号秩检验的关键值制作成表, 而当 $N > 20$ 时, 可以使用正态分布来近似。

19.13.2 比较多个算法

Kruskal-Wallis 检验 (Kruskal-Wallis test) 是 ANOVA 的非参数版本, 是秩检验的多样本推广。给定 $M = L \times N$ 个观测, 例如给定 L 个算法在 N 个数据集上的错误率 $X_{ij} (i = 1, \dots, L, j = 1, \dots, N)$, 我们将它们从小到大排序, 并赋予它们 $1 \sim M$ 的秩 R_{ij} ; 出现平局时还是取平均值。如果原假设

$$H_0: \mu_1 = \mu_2 = \dots = \mu_L$$

成立, 则算法 i 的秩的平均值大约在 1 与 M 的中间, 即 $(M+1)/2$ 。我们将算法 i 样本平均秩记做 $\bar{R}_{i\cdot}$, 并且在平均秩看上去远离中间时拒绝该假设。检验统计量

$$H = \frac{12}{(M+1)L} \sum_{i=1}^L \left(\bar{R}_{i\cdot} - \frac{M+1}{2} \right)^2$$

是自由度为 $L-1$ 的近似卡方分布, 并且当该统计量的值超过 $\chi_{\alpha, L-1}$ 时, 我们拒绝原假设。

与参数方差检验一样, 如果原假设被拒绝, 则我们可以做事后检验, 对秩的逐对比较进行检查。一种做这种检查的方法是 Tukey 检验 (tukey test), 它使用研究者的范围统计量 (studentized range statistic)

$$q = \frac{\bar{R}_{\max} - \bar{R}_{\min}}{\sigma_w}$$

其中, \bar{R}_{\max} 和 \bar{R}_{\min} 分别是 L 个 (秩的) 均值中的最大和最小均值, 而 σ_w^2 是组秩平均值附近的秩的平均方差。我们拒绝组 i 与组 j 具有相同的秩, 接受它们具有不同秩的备择假设, 如果

$$(\bar{R}_{i\cdot} - \bar{R}_{j\cdot}) > q_{\alpha}(L, L(K-1)) \sigma_w$$

其中 $q_{\alpha}(L, L(K-1))$ 被制表。也可以定义单边检验, 以按平均秩对算法排序。

Demsar(2006) 提出使用 CD (critical difference, 关键差) 图进行可视化。在 $1 \sim L$ 的刻度上, 我们标记平均值 $\bar{R}_{i\cdot}$, 并在组之间绘制长度由关键差 $q_{\alpha}(L, L(K-1)) \sigma_w$ 给定的直线, 使得直线连接差不统计显著的组。

19.14 注释

涉及实验设计的材料采用 Montgomery 2005 的讨论, 这里改写成适合机器学习的形式。关于区间估计、假设检验和方差分析的更为详细的讨论可以在任何统计学导论书籍中找到, 如 Ross 1987。

Dietterich(1998) 讨论了各种统计检验方法, 并在多个应用上使用不同的分类算法对其进行了比较。Fawcett(2006) 给出了 ROC 使用和 AUC 计算的综述。Demsar(2006) 给出了在多个数据集上比较分类方法的统计检验综述。

当我们比较两个或多个算法时, 如果具有相同错误率的原假设未被拒绝, 则我们选用最

简单的,即空间和时间复杂度最小的算法。也就是说,如果数据在错误率方面并不偏好任何一个分类算法的时候,则我们使用我们的先验偏好。例如,如果我们对一个线性模型和一个非线性模型进行比较,并且检验没有拒绝二者具有相同的期望错误率,则我们将选择更为简单的线性模型。即使检验拒绝了这样的假设,在选择算法的时候,错误率也仅仅是一个标准。其他标准,如训练的(空间或时间)复杂度、检验的复杂度和可解释性,在实际应用中都可能是更重要的标准。

Yildiz 和 Alpaydin(2006)给出了如何在 MultiTest 算法中使用事后检验结果产生全序。我们做 $L(L-1)/2$ 次单边逐对检验来对 L 个算法定序,但是这些检验很可能不产生全序,而只产生偏序。缺失的链可以使用先验复杂度信息来补充,以便得到全序。使用误差和复杂度这两类信息,拓扑排序可以产生算法的序。

还有一些检验可以检查对比(contrast)。假设 1 和 2 是神经网络方法,而 3 和 4 是模糊逻辑方法。我们可以检验 1 和 2 的平均情况是否不同于 3 和 4 的平均情况,因此可以更一般地比较诸方法。

另一个需要注意的要点是,我们只对误分类率进行评估或比较。这意味着从我们的观点出发,所有的误分类都具有相同的代价。如果事实并非如此,则我们的检验应当基于风险,将一个合适的损失函数考虑在内。这一方面的工作还不是很多。类似地,这些检验也应当从分类推广到回归,使得可以对回归算法的均方误差进行评估,或可以对两个回归算法的误差进行比较。

在比较两个分类算法时,请注意我们只是对它们是否具有相同的期望错误率进行检验。如果是,这也不意味着它们产生相同的错误。这是我们在第 17 章使用的想法:如果不同的分类器产生不同的错误,则我们可以通过组合多个模型来提高准确率。

19.15 习题

1. 在一个两类问题中,假设我们有损失矩阵,其中 $\lambda_{11} = \lambda_{22} = 0$, $\lambda_{21} = 1$, 而 $\lambda_{12} = \alpha$ 。作为 α 的函数,确定决策阈值。
2. 我们可以通过从一个伯努利分布中抽取样本来模拟一个错误概率为 p 的分类器。进行此模拟,并对 $p_0 \in (0, 1)$ 进行二项检验、近似正态检验和 t 检验。对不同的 p 值,将这些检验进行至少 1 000 次并计算拒绝原假设的概率。当 $p_0 = p$ 时,你认为拒绝的期望概率是多少?
3. 假设 $x^i \sim \mathcal{N}(\mu, \sigma^2)$, 其中 σ^2 已知。对假设 $H_0: \mu \geq \mu_0$ 和 $H_1: \mu < \mu_0$ 如何进行检验?
4. K -折交叉验证 t 检验只对错误率的相等性进行检验。如果假设被拒绝,我们并不知道哪个分类算法具有更低的错误率。我们如何对第一个分类算法的错误率不比第二个分类算法更高的假设进行检验?提示:需要对 $H_0: \mu \leq 0$ 和 $H_1: \mu > 0$ 进行检验。
5. 证明总平方和可以分解成组间平方和和组内平方和: $SS_T = SS_b + SS_w$ 。
6. 对符号检验,使用正态分布近似二项分布。
7. 假设有 3 个分类算法。如何将其从最好到最差进行排序?
8. 如果我们有算法 A 的两个变种,算法 B 的三个变种。考虑它们的所有变种,如何比较算法 A 和算法 B 的总体准确率?

9. 提出一种合适的检验, 比较两个回归算法的误差。
10. 提出一种合适的检验, 比较两个增强学习算法的期望奖励。

19.16 参考文献

- Alpaydin, E. 1999. "Combined 5×2 cv F Test for Comparing Supervised Classification Learning Algorithms." *Neural Computation* 11: 1885-1892.
- Bouckaert, R. R. 2003. "Choosing between Two Learning Algorithms based on Calibrated Tests." In *Twentieth International Conference on Machine Learning*, ed. T. Fawcett and N. Mishra, 51-58. Menlo Park, CA: AAAI Press.
- Demsar, J. 2006. "Statistical Comparison of Classifiers over Multiple Data Sets." *Journal of Machine Learning Research* 7: 1-30.
- Dietterich, T. G. 1998. "Approximate Statistical Tests for Comparing Supervised Classification Learning Algorithms." *Neural Computation* 10: 1895-1923.
- Fawcett, T. 2006. "An Introduction to ROC Analysis." *Pattern Recognition Letters* 27: 861-874.
- Montgomery, D. C. 2005. *Design and Analysis of Experiments*. 6th ed., New York: Wiley.
- Ross, S. M. 1987. *Introduction to Probability and Statistics for Engineers and Scientists*. New York: Wiley.
- Turney, P. 2000. "Types of Cost in Inductive Concept Learning." Paper presented at Workshop on Cost-Sensitive Learning at the Seventeenth International Conference on Machine Learning, Stanford University, Stanford, CA, July 2.
- Wolpert, D. H. 1995. "The Relationship between PAC, the Statistical Physics Framework, the Bayesian Framework, and the VC Framework." In *The Mathematics of Generalization*, ed. D. H. Wolpert, 117-214. Reading, MA: Addison-Wesley.
- Yildiz, O. T., and E. Alpaydin. 2006. "Ordering and Finding the Best of $K > 2$ Supervised Learning Algorithms." *IEEE Transactions on Pattern Analysis and Machine Intelligence* 28: 392-402.

概 率 论

我们简略回顾概率论原理、随机变量概念和实例分布。

A. 1 概率论原理

随机试验是其结果不能提前以确定的方式预测的试验 (Ross 1987; Casella 和 Berger 1990)。所有可能的结果的集合称作样本空间 S 。一个样本空间是离散的, 如果它由结果的有限(或可数无限)集组成; 否则是连续的。 S 的任意子集 E 是一个事件。事件是集合, 并且我们可以谈论它们的补、交、并等。

概率的一种解释是频率 (frequency)。当一个试验在完全相同的条件下不断重复时, 对于任意事件 E , 结果在 E 中的次数所占的比例趋向于某个常数值。这个常数极限频率是事件的概率, 而我们把它记作 $P(E)$ 。

有时, 概率可解释成可信程度 (degree of belief)。例如, 当我们说土耳其赢得 2010 年足球世界杯冠军的概率时, 我们并不是指出现的频率, 因为 2010 年足球世界杯只进行一次, 并且(在写本书时)它还未进行。在这种情况下, 我们的意思是我们主观相信该事件出现的程度。由于是主观的, 因此对同一事件, 不同的人可能指派不同的概率。

517

A. 1. 1 概率论公理

公理确保随机试验中指派的概率可以解释成相对频率, 并且这些指派符合我们对相对频率之间关系的直观理解:

1) $0 \leq P(E) \leq 1$ 。如果 E_1 是不可能出现的事件, 则 $P(E_1) = 0$ 。如果 E_2 是一定出现的事件, 则 $P(E_2) = 1$ 。

2) 如果 S 是包含所有可能结果的样本空间, 则 $P(S) = 1$ 。

3) 如果 $E_i, i=1, \dots, n$, 是互斥的(即如果它们不可能同时出现: $E_i \cap E_j = \emptyset, i \neq j$, 其中 \emptyset 是不包含任何可能结果的空事件), 则我们有

$$P\left(\bigcup_{i=1}^n E_i\right) = \sum_{i=1}^n P(E_i) \quad (\text{A. 1})$$

例如, 设 E^c 表示 E 的补, 由不在 E 中的 S 中所有可能的结果组成, 我们有 $E \cap E^c = \emptyset$, 并且

$$\begin{aligned} P(E \cup E^c) &= P(E) + P(E^c) = 1 \\ P(E^c) &= 1 - P(E) \end{aligned}$$

如果 E 和 F 的交非空, 则我们有

$$P(E \cup F) = P(E) + P(F) - P(E \cap F) \quad (\text{A. 2})$$

A. 1.2 条件概率

$P(E|F)$ 是给定事件 F 出现的情况下, 事件 E 出现的概率, 并由下式给出

$$P(E|F) = \frac{P(E \cap F)}{P(F)} \quad (\text{A. 3})$$

知道事件 F 出现将样本空间缩小到 F , 而 E 也出现的部分为 $E \cap F$ 。注意, 式(A.3)仅当 $P(F) > 0$ 时才有定义。由于 \cap 是可交换的, 我们有

$$P(E \cap F) = P(E|F)P(F) = P(F|E)P(E) \quad 518$$

由此得到贝叶斯公式 (Bayes' formula):

$$P(F|E) = \frac{P(E|F)P(F)}{P(E)} \quad (\text{A. 4})$$

当 F_i 互斥并穷举时, 即当 $\bigcup_{i=1}^n F_i = S$ 时

$$E = \bigcup_{i=1}^n E \cap F_i$$

$$P(E) = \sum_{i=1}^n P(E \cap F_i) = \sum_{i=1}^n P(E|F_i)P(F_i) \quad (\text{A. 5})$$

贝叶斯公式使得我们可以有

$$P(F_i|E) = \frac{P(E \cap F_i)}{P(E)} = \frac{P(E|F_i)P(F_i)}{\sum_j P(E|F_j)P(F_j)} \quad (\text{A. 6})$$

如果 E 和 F 是独立的 (independent), 则我们有 $P(E|F) = P(E)$, 因此

$$P(E \cap F) = P(E)P(F) \quad (\text{A. 7})$$

也就是说, F 是否出现并不改变 E 出现的概率。

A. 2 随机变量

随机变量 (random variable) 是一个函数, 它为随机试验的样本空间中的每个结果指派一个数。

A. 2.1 概率分布与密度函数

对于任意实数值 a , 随机变量 X 的概率分布函数 (probability distribution function) $F(\cdot)$ 是

$$F(a) = P\{X \leq a\} \quad (\text{A. 8})$$

并且我们有

$$P\{a < X \leq b\} = F(b) - F(a) \quad (\text{A. 9})$$

如果 X 是离散的随机变量, 则

$$F(a) = \sum_{x \leq a} P(x) \quad (\text{A. 10}) \quad 519$$

其中 $P(\cdot)$ 是概率质量函数 (probability mass function), 定义为 $P(a) = P\{X = a\}$ 。如果 X 是连续的随机变量, 则 $p(\cdot)$ 是概率密度函数 (probability density function), 使得

$$F(a) = \int_{-\infty}^a p(x) dx \quad (\text{A. 11})$$

A. 2.2 联合分布与密度函数

在特定的试验中，我们可能对两个或多个随机变量之间的关系感兴趣，并且我们使用 X 和 Y 的联合 (joint) 概率分布和密度函数，满足

$$F(x, y) = P\{X \leq x, Y \leq y\} \quad (\text{A. 12})$$

单个边缘 (marginal) 分布和密度可以通过边缘化来计算，即在自由变量上求和：

$$F_X(x) = P\{X \leq x\} = P\{X \leq x, Y \leq \infty\} = F(x, \infty) \quad (\text{A. 13})$$

在离散情况下，我们有

$$P(X = x) = \sum_j P(x, y_j) \quad (\text{A. 14})$$

而在连续情况下，我们有

$$p_X(x) = \int_{-\infty}^{\infty} p(x, y) dy \quad (\text{A. 15})$$

如果 X 和 Y 是独立的 (independent)，我们有

$$p(x, y) = p_X(x)p_Y(y) \quad (\text{A. 16})$$

这些都能够以直截了当的方式推广到多于两个随机变量的情况。

A. 2.3 条件分布

当 X 和 Y 是随机变量时，

$$P_{X|Y}(x|y) = P\{X = x | Y = y\} = \frac{P\{X = x, Y = y\}}{P\{Y = y\}} = \frac{P(x, y)}{P_Y(y)} \quad (\text{A. 17})$$

A. 2.4 贝叶斯规则

当两个随机变量联合分布，其中一个的值已知时，另一个取给定值的概率可以使用贝叶斯规则计算：

$$P(y|x) = \frac{P(x|y)P_Y(y)}{P_X(x)} = \frac{P(x|y)P_Y(y)}{\sum_y P(x|y)P_Y(y)} \quad (\text{A. 18})$$

换言之，

$$\text{后验} = \frac{\text{似然} \times \text{先验}}{\text{证据}} \quad (\text{A. 19})$$

注意，分母通过在所有可能的 y 值上对分子求和 (或积分，如果 y 是连续的) 得到。 $p(y|x)$ 的“形状”取决于分子，分母作为规范化因子确保 $p(y|x)$ 的和为 1。通过考虑 x 提供的信息，贝叶斯规则使得我们将一个先验概率修改为后验概率。

贝叶斯规则反转依赖性，如果 $p(x|y)$ 已知，使得我们可以计算 $p(y|x)$ 。假设 y 是 x 的“原因”，如 y 是度暑假， x 是被晒黑，则 $p(x|y)$ 是已知某人度暑假，他被晒黑的概率。这是

因果(causal)(或预测)方法。贝叶斯规则允许我们使用诊断(diagnostic)方法来计算 $p(y|x)$: 即某人知道会被晒黑,他去度暑假的概率。 $p(y)$ 是任何人去度暑假的概率,而 $p(x)$ 是任何人被晒黑的概率,包括度暑假和不度暑假的人。

A. 2.5 期望

随机变量 X 的期望(expectation)、期望值(expected value)或均值(mean)记作 $E[X]$,是大量试验中 X 的平均值:

$$E[X] = \begin{cases} \sum_i x_i P(x_i) & \text{如果 } X \text{ 是离散的} \\ \int x p(x) dx & \text{如果 } X \text{ 是连续的} \end{cases} \quad (\text{A. 20})$$

它是加权平均,其中每个值被 X 取该值的概率加权。它具有如下性质($a, b \in \mathfrak{R}$):

$$\begin{aligned} E[aX + b] &= aE[X] + b \\ E[X + Y] &= E[X] + E[Y] \end{aligned} \quad (\text{A. 21}) \quad \boxed{521}$$

对于任意实数值函数 $g(\cdot)$,期望值是

$$E[g(X)] = \begin{cases} \sum_i g(x_i) P(x_i) & \text{如果 } X \text{ 是离散的} \\ \int g(x) p(x) dx & \text{如果 } X \text{ 是连续的} \end{cases} \quad (\text{A. 22})$$

一种特例 $g(x) = x^n$,称作 X 的 n 阶矩,定义为

$$E[X^n] = \begin{cases} \sum_i x_i^n P(x_i) & \text{如果 } X \text{ 是离散的} \\ \int x^n p(x) dx & \text{如果 } X \text{ 是连续的} \end{cases} \quad (\text{A. 23})$$

均值(mean)是一阶矩并记作 μ 。

A. 2.6 方差

方差(variance)度量 X 在期望值附近的变化。如果 $\mu \equiv E[X]$,则方差定义为

$$\text{Var}(X) = E[(X - \mu)^2] = E[X^2] - \mu^2 \quad (\text{A. 24})$$

方差是二阶矩减去一阶矩的平方。方差记作 σ^2 ,具有如下性质($a, b \in \mathfrak{R}$):

$$\text{Var}(aX + b) = a^2 \text{Var}(X) \quad (\text{A. 25})$$

$\sqrt{\text{Var}(X)}$ 称作标准差(standard deviation),记作 σ 。标准差具有和 X 相同的单位,并且比方差容易解释。

协方差(covariance)指示两个随机变量之间的关系。如果 X 的出现使得 Y 更可能出现,则协方差为正;如果 X 的出现使得 Y 更不可能发生,则协方差为负;如果没有依赖性,则协方差为0。

$$\text{Cov}(X, Y) = E[(X - \mu_X)(Y - \mu_Y)] = E[XY] - \mu_X \mu_Y \quad (\text{A. 26})$$

其中 $\mu_X \equiv E[X]$, $\mu_Y \equiv E[Y]$ 。一些其他性质是

$$\text{Cov}(X, Y) = \text{Cov}(Y, X)$$

$$\text{Cov}(X, X) = \text{Var}(X)$$

$$\text{Cov}(X+Z, Y) = \text{Cov}(X, Y) + \text{Cov}(Z, Y)$$

$$\text{Cov}\left(\sum_i X_i, Y\right) = \sum_i \text{Cov}(X_i, Y) \quad (\text{A. 27})$$

$$\text{Var}(X+Y) = \text{Var}(X) + \text{Var}(Y) + 2\text{Cov}(X, Y) \quad (\text{A. 28})$$

$$\text{Var}\left(\sum_i X_i\right) = \sum_i \text{Var}(X_i) + \sum_i \sum_{j \neq i} \text{Cov}(X_i, X_j) \quad (\text{A. 29})$$

如果 X 和 Y 是独立的, 则 $E[XY] = E[X]E[Y] = \mu_X \mu_Y$, 并且 $\text{Cov}(X, Y) = 0$ 。这样, 如果 X_i 是独立的, 则

$$\text{Var}\left(\sum_i X_i\right) = \sum_i \text{Var}(X_i) \quad (\text{A. 30})$$

相关性(correlation)是一个规范化的、维无关的量, 其值总是在 -1 和 1 之间:

$$\text{Corr}(X, Y) = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X)\text{Var}(Y)}} \quad (\text{A. 31})$$

A. 2.7 弱大数定律

设 $X = \{X^t\}_{t=1}^N$ 是独立的同分布的(iid)随机变量的集合, 每个都具有均值 μ 和有限方差 σ^2 。则对于任意 $\varepsilon > 0$,

$$P\left\{\left|\frac{\sum_t X^t}{N} - \mu\right| > \varepsilon\right\} \rightarrow 0 \quad \text{随 } N \rightarrow \infty \quad (\text{A. 32})$$

也就是说, 随着 N 趋向于无穷大, N 个试验的平均值趋向于均值。

A. 3 特殊的随机变量

有一些类型的随机变量频繁出现, 因此对它们命名。

A. 3.1 伯努利分布

试验进行, 其结果或者“成功”, 或者“失败”。随机变量 X 是一个 0/1 指示变量, 并且对于成功结果取值 1, 否则为 0。 p 是试验结果为成功的概率。则

$$P\{X = 1\} = p, \quad \text{而 } P\{X = 0\} = 1 - p \quad (\text{A. 33})$$

这等价于

$$P\{X = i\} = p^i(1-p)^{1-i}, i = 0, 1 \quad (\text{A. 34})$$

如果 X 是伯努利变量, 则它的期望值和方差是

$$E[X] = p, \quad \text{Var}(X) = p(1-p) \quad (\text{A. 35})$$

A. 3.2 二项分布

如果做了 N 次相同的、独立的伯努利试验, 代表 N 次试验中成功次数的随机变量 X 是二项分布的。 i 次成功的概率为

$$P\{X = i\} = C_N^i p^i (1-p)^{N-i}, i = 0 \cdots N \quad (\text{A. 36})$$

如果 X 是二项的, 则它的期望值和方差为

$$E[X] = Np, \quad \text{Var}(X) = Np(1-p) \quad (\text{A. 37})$$

A. 3.3 多项分布

考虑伯努利分布的推广。其中, 取代两种状态, 随机事件的结果是 K 个互斥、穷举状态之一, 每个具有出现概率 p_i , 其中 $\sum_{i=1}^K p_i = 1$ 。假设做了 N 次这样的试验, 其中结果 i 出现 N_i 次, 满足 $\sum_{i=1}^K N_i = N$ 。则 N_1, N_2, \dots, N_K 的联合分布是多项分布:

$$P(N_1, N_2, \dots, N_K) = N! \prod_{i=1}^K \frac{p_i^{N_i}}{N_i!} \quad (\text{A. 38})$$

当 $N=1$ 时是一种特殊情况: 只做了一次试验。于是 N_i 是 0/1 指示变量, 其中只有一个为 1, 其余均为 0。式(A. 38)归约为

$$P(N_1, N_2, \dots, N_K) = \prod_{i=1}^K p_i^{N_i} \quad (\text{A. 39})$$

A. 3.4 均匀分布

X 均匀地分布在区间 $[a, b]$ 上, 如果它的密度函数由下式给定

$$p(x) = \begin{cases} \frac{1}{b-a} & \text{如果 } a \leq x \leq b \\ 0 & \text{否则} \end{cases} \quad (\text{A. 40})$$

如果 X 是均匀的, 则它的期望值和方差为

$$E[X] = \frac{a+b}{2}, \quad \text{Var}(X) = \frac{(b-a)^2}{12} \quad (\text{A. 41}) \quad \boxed{524}$$

A. 3.5 正态(高斯)分布

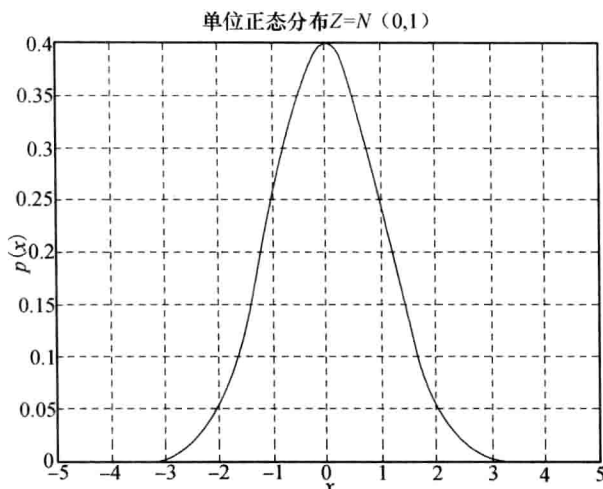
X 是均值为 μ 、方差为 σ^2 的正态或高斯分布, 记作 $\mathcal{N}(\mu, \sigma^2)$, 如果它的密度函数是

$$p(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{(x-\mu)^2}{2\sigma^2}\right], \quad -\infty < x < \infty \quad (\text{A. 42})$$

许多随机现象都遵守钟形正态分布, 或至少近似地遵守正态分布; 许多自然观测都可以看作连续的、典型值的稍微不同的版本——这或许是将它称作正态(normal)分布的原因。在这种情况下, μ 定义典型值, 而 σ 定义典型值附近实例变化的大小。

68.27% 的值落在 $(\mu - \sigma, \mu + \sigma)$ 中, 95.45% 的值落在 $(\mu - 2\sigma, \mu + 2\sigma)$ 中, 99.73% 的值落在 $(\mu - 3\sigma, \mu + 3\sigma)$ 中。这样, $P\{|x - \mu| < 3\sigma\} \approx 0.99$ 。实践中, 如果 $x < \mu - 3\sigma$ 或 $x > \mu + 3\sigma$, 则 $p(x) \approx 0$ 。 Z 是单位正态分布, 即 $\mathcal{N}(0, 1)$ (见图 A-1), 并且它的密度记作

$$p_Z(x) = \frac{1}{\sqrt{2\pi}} \exp\left[-\frac{x^2}{2}\right] \quad (\text{A. 43}) \quad \boxed{525}$$

图 A-1 单位正态分布 Z 的概率密度函数

如果 $X \sim \mathcal{N}(\mu, \sigma^2)$ 并且 $Y = aX + b$, 则 $Y \sim \mathcal{N}(a\mu + b, a^2\sigma^2)$ 。独立的正态变量的和也是正态的, 其中 $\mu = \sum_i \mu_i$, $\sigma^2 = \sum_i \sigma_i^2$ 。如果 X 是 $\mathcal{N}(\mu, \sigma^2)$, 则

$$\frac{X - \mu}{\sigma} \sim Z \quad (\text{A. 44})$$

这称作 z -标准化。

设 X_1, X_2, \dots, X_N 是 iid 随机变量, 都具有均值 μ 和方差 σ^2 。则中心极限定理 (central limit theorem) 表明对于大的 N , 分布

$$X_1 + X_2 + \dots + X_N \quad (\text{A. 45})$$

近似于 $\mathcal{N}(N\mu, N\sigma^2)$ 。例如, 如果 X 是参数为 (N, p) 的二项分布, 则 X 可以写成 N 个伯努利试验的和, 并且 $(X - Np) / \sqrt{Np(1-p)}$ 是近似单位正态的。

中心极限定理也用来在计算机上产生正态分布的随机变量。程序设计语言具有一些子程序, 返回 $[0, 1]$ 上均匀分布的 (伪) 随机数。当 U_i 是这样的随机变量时, $\sum_{i=1}^{12} U_i - 6$ 近似于 Z 。

设 $X^t \sim \mathcal{N}(\mu, \sigma^2)$ 。估计样本均值

$$m = \frac{\sum_{t=1}^N X^t}{N} \quad (\text{A. 46})$$

也是正态的, 均值为 μ , 而方差为 σ^2/N 。

A. 3.6 卡方分布

如果 Z_i 是独立的单位正态随机变量, 则

$$X = Z_1^2 + Z_2^2 + \dots + Z_n^2 \quad (\text{A. 47})$$

是自由度为 n 的卡方分布, 即 $X \sim \chi_n^2$, 其中

$$E[X] = n, \text{Var}(X) = 2n \quad (\text{A. 48})$$

当 $X' \sim \mathcal{N}(\mu, \sigma^2)$ 时, 估计样本方差为

$$S^2 = \frac{\sum_i (X' - m)^2}{N - 1} \quad (\text{A. 49})$$

并且我们有

$$(N - 1) \frac{S^2}{\sigma^2} \sim \chi_{N-1}^2 \quad (\text{A. 50})$$

还知道 m 和 S^2 是独立的。

526

A. 3.7 t 分布

如果 $Z \sim \mathcal{N}$ 和 $X \sim \chi_n^2$ 是独立的, 则

$$T_n = \frac{Z}{\sqrt{X/n}} \quad (\text{A. 51})$$

是自由度为 n 的 t -分布, 其中

$$E[T_n] = 0, n > 1, \quad \text{Var}(T_n) = \frac{n}{n-2}, n > 2 \quad (\text{A. 52})$$

像单位正态密度一样, t -分布在 0 周围是对称的。随着 n 越来越大, t 密度变得越来越像正态分布, 区别是 t -分布具有较粗的尾部, 表明比正态分布具有更大的可变性。

A. 3.8 F 分布

如果 $X_1 \sim \chi_n^2$ 和 $X_2 \sim \chi_m^2$ 分别是自由度为 n 和 m 的卡方随机变量, 则

$$F_{n,m} = \frac{X_1/n}{X_2/m} \quad (\text{A. 53})$$

是自由度为 n 和 m 的 F -分布, 其中

$$E[F_{n,m}] = \frac{m}{m-2}, m > 2, \quad \text{Var}(F_{n,m}) = \frac{m^2(2m+2n-4)}{n(m-2)^2(m-4)}, m > 4 \quad (\text{A. 54})$$

A. 4 参考文献

Casella, G., and R. L. Berger. 1990. *Statistical Inference*. Belmont, CA: Duxbury.

Ross, S. M. 1987. *Introduction to Probability and Statistics for Engineers and Scientists*. New York: Wiley.

527

索引

索引中标注的页码为英文原书页码, 与书中边栏的页码一致。

0/1 loss function(0/1 损失函数), 51

5×2

cross-validation(~ 交叉验证), 488

cv paired F test(~ 交叉验证与配对 F 检验), 503

cv paired t test(~ 交叉验证与配对 t 检验), 503

A

Active learning(主动学习), 360

AdaBoost, 431

Adaptive resonance theory(自适应共鸣理论), 285

Additive models(加法模型), 180

Agglomerative clustering(凝聚聚类), 157

AIC, *see* Akaike's information criterion 参见 Akaike 信息准则

Akaike's information criterion(Akaike 信息准则), 81

Alignment(比对), 324

Analysis of variance(方差分析), 504

Anchor(锚), 291

ANOVA, *see* Analysis of variance 参见方差分析

Approximate normal test(近似正态检验), 500

Apriori algorithm(Apriori 算法)56

Area under the curve(曲线下方面积), 491

ART, *see* Adaptive resonance theory 参见自适应共鸣理论

Artificial neural networks(人工神经网络), 233

Association rule(关联规则), 4, 55

Attribute(属性), 87

Autoassociator(自动关联器), 268

B

Backpropagation(后向传播), 250

through time(通过时间 ~), 272

Backup(后退), 456

Backward selection(向后选择), 111

Backward variable(向后变量), 372

Bag of words(词袋), 102, 324

Bagging(装袋), 430

Base-learner(基学习器), 419

Basis function(基函数), 211

cooperative vs. competitive(协作与竞争 ~), 297

for a kernel(内核), 352

normalization(规范化 ~), 295

Basket analysis(购物篮分析), 55

Batch learning(批学习), 251

Baum-Welch algorithm(Baum-Welch 算法), 376

Bayes' ball(贝叶斯球), 402

Bayes' classifier(贝叶斯分类(器)), 51

Bayes' estimator(贝叶斯估计(器)), 68

Bayes' rule(贝叶斯规则), 49, 521

Bayesian information criterion(贝叶斯信息准则), 81

Bayesian model combination(贝叶斯模型组合), 426

Bayesian model selection(贝叶斯模型选择), 82

Bayesian network(贝叶斯网络), 387

Belief networks(信念网络), 387

Belief state(信任状态), 465

Bellman's equation(Bellman 公式), 452

Beta distribution(贝塔分布), 345

Between-class scatter matrix(类间散布矩阵), 130

Bias(偏倚), 65

Bias unit(偏倚单元), 238

Bias/variance dilemma(偏倚/方差两难), 78

BIC, *see* Bayesian information criterion 参见贝叶斯信息准则

Binary split(二元划分), 187

Binding(绑定), 202

Binomial test(二项检验), 499

Biometrics(生物测定学), 441

Blocking(阻塞), 482

Bonferroni correction(Bonferroni 校正), 508

Boosting(提升), 431
 Bootstrap(自助法), 489

C

C4.5, 191
 C4.5Rules(C4.5 规则), 197
 CART, 191, 203
 Cascade correlation(级联相关), 264
 Cascading(级联), 438
 Case-based reasoning(基于案例的推理), 180
 Causality(因果关系), 396
 causal graph(因果图), 388
 Central limit theorem(中心极限定理), 526
 Class(类)
 confusion matrix(~ 混淆矩阵), 493
 likelihood(~ 似然), 50
 Classification(分类), 5
 likelihood- vs. discriminant-based(基于似然与基于判别式的 ~), 209
 Classification tree(分类树), 188
 Clique(团), 411
 Cluster(簇), 144
 Clustering(聚类), 11
 agglomerative(凝聚 ~), 157
 divisive(分裂 ~), 157
 hierarchical(层次 ~), 157
 online(在线 ~), 281
 Code word(码字), 146
 Codebook vector(编码本向量), 146
 Coefficient of determination (of regression)(回归判别式的系数), 76
 Color quantization(颜色量化), 145
 Common principal components(公共主成分), 119
 Competitive basis functions(竞争的基函数), 297
 Competitive learning(竞争学习), 280
 Complete-link clustering(全链接聚类), 158
 Component density(支密度), 144
 Compression(压缩), 8, 146
 Condensed nearest neighbor(精简的最近邻), 173
 Confidence interval(置信区间)
 one-sided(单侧 ~), 495
 two-sided(双侧 ~), 494

Confidence of an association rule(关联规则的置信度), 55
 Conjugate prior(共轭先验), 344
 Connection weight(连接权重), 237
 Contingency table(列联表), 501
 Correlation(相关), 89
 Cost-sensitive learning(代价敏感学习), 478
 Coupled HMM(耦合的 HMM)400
 Covariance function(协方差函数)358
 Covariance matrix(协方差矩阵), 88
 Credit assignment(信用分配), 448
 Critic(批评者), 447
 CRM, *see* Customer relationship management 参见客户关系管理
 Cross-entropy(互熵), 221
 Cross-validation(交叉验证), 40, 80, 486
 $5 \times 2(5 \times 2 \sim)$, 488
 K-fold(K-折 ~), 487
 Curse of dimensionality(维(度)灾难), 160
 Customer relationship management(客户关系管理), 155
 Customer segmentation(客户段), 155

D

d-seperation(d-分离), 402
 Decision node(决策节点), 185
 Decision region(决策区域), 53
 Decision tree(决策树), 185
 multivariate, (多元 ~)202
 omnivariate(杂变量 ~), 205
 soft(软 ~), 305
 univariate(单变量 ~), 187
 Delve repository(Delve 知识库), 17
 Dendrogram(系统树图), 158
 Density estimation(密度估计), 11
 Dichotomizer(两分器), 53
 Diffusion kernel(扩散核), 325
 Dimensionality reduction(维度归约)
 nonlinear(非线性 ~), 269
 Directed acyclic graph(有向无环图), 387
 Dirichlet distribution(狄里克雷分布), 344
 Discount rate(折扣率), 451
 Discriminant(判别式), 5

function(~ 函数), 53
 linear(线性 ~), 97
 quadratic(二次 ~), 95
 Discriminant adaptive nearest neighbor(判别式自适应最近邻), 172
 Discriminant-based classification(基于判别式的分类), 209
 Distributed vs. local representation(分布的与局部的表示), 156, 287
 Diversity(多样性)420
 Divisive clustering(分裂聚类), 157
 Document categorization(文档分类), 102
 Doubt(不确定(实例)), 26
 Dual representation(对偶表示), 337, 352
 Dynamic classifier selection(动态分类器选择), 435
 Dynamic node creation(动态节点创建), 264
 Dynamic programming(动态规划), 453

E

Early stopping(提前终止), 223, 258
 ECOC, *see* Error-correcting output codes 参见纠错输出码
 Edit distance(编辑距离), 324
 Eigendigits(特征数字), 118
 Eigenfaces(特征面孔), 118
 Eligibility trace(资格迹), 459
 EM, *see* Expectation-Maximization 参见期望最大化
 Emission probability(发射概率), 367
 Empirical error(经验误差), 24
 Ensemble(系综), 424
 Entropy(熵), 188
 Episode(片段), 451
 Epoch(周期), 251
 Error(误差, 错误)
 type I(第一类 ~), 497
 type II(第二类 ~), 497
 Error-correcting output codes(纠错输出码), 427
 Euclidean distance(欧氏距离), 98
 Evidence(证据), 50
 Example(实例), 87
 Expectation-Maximization(期望最大化), 150
 supervised(监督的 ~), 299

Expected error rate(期望误差率), 476
 Expected utility(期望效用), 54
 Experiment(实验)
 design(~ 设计), 478
 factorial(~ 因素)481
 strategies(~ 策略)480
 Explaining away(解释远离), 50
 Extrapolation(外推), 29

F

FA, *see* Factor analysis 参见因素分析
 Factor analysis(因素分析), 120
 Factor graph(因素图), 412
 Factorial HMM(因素 HMM)400
 Feature(特征), 87
 extraction(~ 提取), 110
 selection(~ 选择), 110
 Finite-horizon(有限视野), 451
 First-order rule(一阶规则), 201
 Fisher's linear discriminant(费希尔线性判别式), 129
 Flexible discriminant analysis(柔性判别式分析), 120
 Floating search(浮动搜索), 112
 Foil(Foil 算法), 199
 Forward selection(向前选择), 110
 Forward variable(正向变量), 370
 Forward-backward procedure(正反向过程), 370
 Fuzzy k -means(模糊 k -均值), 160
 Fuzzy membership function(模糊隶属函数), 295
 Fuzzy rule(模糊规则), 295

G

Gamma distribution(伽马分布), 347
 Gamma function(伽马函数)344
 Gaussian prior(高斯先验), 349
 Generalization(推广, 泛化), 24, 39
 Generalized linear models(广义线性模型), 230
 Generative model(生成模型)342, 397
 Generative topographic mapping(生成拓扑映射), 306
 Geodesic distance(测地距离), 133
 Gini index(基尼指数), 189
 Gradient descent(梯度下降), 219
 stochastic(随机 ~), 241

Gradient vector(梯度向量), 219
 Gram matrix(Gram 矩阵), 321
 Graphical models(图模型), 387
 Group(分组), 144
 GTM, *see* Generative topographic mapping 参见生成
 拓扑映射

H

Hamming distance(汉明距离), 171
 Hebbian learning(Hebbian 学习), 283
 Hidden layer(隐藏层), 246
 Hidden Markov model(隐马尔可夫模型), 367, 398
 coupled(耦合 ~), 400
 factorial(因子 ~), 400
 input-output(输入-输出 ~), 379, 400
 left-to-right(自左向右 ~), 380
 switching(开关 ~), 400
 Hidden variables(隐藏变量), 57, 396
 Hierarchical clustering(层次聚类), 157
 Hierarchical cone(层次锥体), 260
 Hierarchical mixture of experts(层次混合专家模型), 304
 Higher-order term(高阶项), 211
 Hinge loss(Hinge 损失), 317
 Hint(线索), 261
 Histogram(直方图), 165
 HMM, *see* Hidden Markov model 参见隐马尔可夫
 模型
 Hybrid learning(混合学习), 291
 Hypothesis(假设), 23
 class(~类), 23
 most general(最一般的 ~), 24
 most specific(最特殊的 ~), 24
 Hypothesis testing(假设检验), 496

I

ID3(ID3 算法), 191
 IF-THEN rules(IF-THEN 规则), 197
 Iid(independent and identically distributed)(独立同
 分布), 41
 Ill-posed(不适定的), 38
 Impurity measure(不纯度度量), 188

Imputation(派算), 89
 Independence(独立性), 388
 Inductive bias(归纳偏倚), 38
 Inductive logic programming(归纳逻辑程序设计), 202
 Infinite-horizon(无限视野), 451
 Influence diagrams(影响图), 414
 Information retrieval(信息检索), 491
 Initial probability(初始概率), 364
 Input(输入), 87
 Input representation(输入表示), 21
 Input-output HMM(输入-输出 HMM), 379, 399
 Instance(实例), 87
 Instance-based learning(基于实例的学习), 164
 Interest of an association rule(关联规则的兴趣
 度), 55
 Interpolation(插值), 35
 Interpretability(可解释性), 197
 Interval estimation(区间估计), 493
 Irep(Irep(算法)), 199
 Isometric feature mapping(等距特征映射), 133

J

Job shop scheduling(作业车间调度), 471
 Junction tree(结树), 410

K

K-armed bandit(K 臂赌博机), 449
 K-fold(K-折)
 cross-validation(~交叉验证), 487
 cv paired *t* test(~交叉验证与配对 *t* 检验), 502
k-means clustering(*k*-均值聚类), 147
 fuzzy(模糊 ~), 172
 online(在线 ~), 281
k-nearest neighbor(*k*-最近邻)
 classifier(~分类), 172
 density estimate(~估计), 169
 smoother(~光滑), 177
k-nn, *see* *k*-nearest neighbor 参见 *k*-最近邻
 Kalman filter(Kalman 过滤), 400
 Karhunen-Loève expansion(Karhunen-Loève 展
 开), 119

Kernel estimator(核估计), 167
 Kernel function(核函数), 167, 320, 353
 Kernel PCA(核PCA), 336
 Kernel smoother(核光滑), 176
 Knowledge extraction(知识提取), 8, 198, 295
 Kolmogorov complexity(Kolmogorov 复杂度), 82
 Kruskal-Wallis test(Kruskal-Wallis 检验), 511

L

Laplace approximation(拉普拉斯近似), 354
 Laplace prior(拉普拉斯先验), 350
 Lasso, 352
 Latent factors(潜在因子), 120
 Lateral inhibition(横向抑制), 282
 LDA, *see* Linear discriminant analysis 参见线性判别分析
 Leader cluster algorithm(领导者聚类算法), 147
 Leaf node(叶节点), 186
 Learning automata(学习自动机), 471
 Learning vector quantization(学习向量量化), 300
 Least squares difference test(最小方差检验), 507
 Leave-one-out(留一), 487
 Least squares estimate(最小方差估计), 74
 Left-to-right HMM(自左向右 HMM), 380
 Level of significance(显著水平), 497
 Levels of analysis(分析层面), 234
 Lift of an association rule(关联规则的提升度), 55
 Likelihood(似然), 62
 Likelihood ratio(似然率), 58
 Likelihood-based classification(基于似然的分类), 209
 Linear classifier(线性分类器), 97, 210
 Linear discriminant(线性判别式), 97, 210
 Linear discriminant analysis(线性判别分析), 128
 Linear dynamical system(线性动态系统), 400
 Linear opinion pool(线性判断组合), 424
 Linear regression(线性回归), 74
 multivariate(多元~), 103
 Linear separability(线性可分性), 215
 Local representation(局部表示), 288
 Locally linear embedding(局部线性嵌入), 135
 Locally weighted running line smoother(局部加权移动

线性光滑), 177
 Loess, *see* Locally weighted running line smoother 参见局部加权移动线性光滑
 Log likelihood(对数似然), 62
 Log odds(对数几率), 58, 218
 Logistic discrimination(逻辑斯谛判别式), 220
 Logistic function(逻辑斯谛函数), 218
 Logit(分对数), 218
 Loss function(损失函数), 51
 LVQ, *see* Learning vector quantization 参见学习向量量化

M

Mahalanobis distance(马氏距离), 90
 Margin(边缘), 25, 311, 433
 Markov decision process(马尔可夫决策过程), 451
 Markov mixture of experts(马尔可夫混合专家模型), 379
 Markov model(马尔可夫模型), 364
 hidden(隐~), 367
 learning(学习~), 366, 375
 observable(可观测的~), 365
 Maximum a posteriori (MAP) estimate(最大化后验估计), 68, 343
 Maximum likelihood estimation(最大似然估计), 62
 McNemar's test(McNemar 检验), 501
 MDP, *see* Markov decision process 参见马尔可夫决策过程
 MDS, *see* Multidimensional scaling 参见多维定标
 Mean square error(均方误差), 65
 Mean vector(均值向量), 88
 Memory-based learning(基于记忆的学习), 164
 Minimum description length(最大描述长度), 82
 Mixture components(混合分支), 144
 Mixture density(混合密度), 144
 Mixture of experts(混合专家模型), 301, 434
 competitive(竞争的~), 304
 cooperative(协作的~), 303
 hierarchical(层次的~), 305
 Markov(马尔可夫~), 379, 400
 Mixture of factor analyzers(混合因子分析方法), 155
 Mixture of mixtures(混合的混合(密度)), 156

Mixture of probabilistic principal component analyzers (混合概率主成分分析), 155

Mixture proportion(混合比例), 144

MLE, *see* Maximum likelihood estimation 参见最大似然估计

Model combination(模型组合)

- multiexpert(多专家~), 423
- multistage(多级~), 423

Model selection(模型选择), 38

MoE, *see* Mixture of experts 参见混合专家模型

Momentum(动量), 257

Moralization(教化), 411

Multidimensional scaling(多维定标), 125

- nonlinear(非线性~), 287
- using MLP(使用 MLP~), 269

Multilayer perceptrons(多层感知器), 246

Multiple comparisons(多重比较), 507

Multivariate linear regression(多元线性回归), 103

Multivariate polynomial regression(多元多项式回归), 104

Multivariate tree(多变量/元树), 202

N

Naive Bayes' classifier(朴素贝叶斯分类), 397

- discrete input(离散输入~), 102
- numeric input(数值输入~), 97

Naive estimator(朴素估计), 166

Nearest mean classifier(最近均值分类(器)), 98

Nearest neighbor classifier(最近邻分类(器)), 166

- condensed(精简的~), 173

Negative examples(负实例), 21

Neuron(神经元), 233

No Free Lunch Theorem(没有免费的午餐法则), 477

Noisy OR(噪声或), 109

Noise(噪声), 30

Nonparametric estimation(非参数估计), 163

Nonparametric testing(非参数检验), 508

Null hypothesis(原假设), 497

O

Observable Markov model(可观测的马尔可夫模

型), 365

Observable variable(可观测变量), 48

Observation(观测), 87

Observation probability(观测概率), 367

OCI(OCI(算法)), 203

Occam's razor(奥克姆剃刀), 32

Off-policy(无策略), 457

Omnivariate decision tree(杂变量决策树), 205

On-policy(有策略), 458

One-sided confidence interval(单侧置信区间), 495

One-sided test(单侧检验), 498

Online k -means(在线 k -均值), 281

Online learning(在线学习), 241

Optimal policy(最优策略), 452

Optimal separating hyperplane(最佳分离超平面), 311

Outlier detection(离群点检测), 9, 333

Overfitting(过(分)拟合), 39, 79

Overtraining(过(分)训练), 258

P

PAC, *see* Probably Approximately Correct 参见或然近似正确

Paired test(配对检验), 501

Pairwise separation(逐对分离), 216, 428

Parallel processing(并行处理), 236

Partially observable Markov decision process(部分可观测的马尔可夫决策过程), 464

Parzen windows(Parzen 窗口), 167

Pattern recognition(模式识别), 6

PCA, *see* Principal components analysis 参见主成分分析

Pedigree(家谱)400

Perceptron(感知器), 237

Phone(音素), 381

Piecewise approximation(分段近似)

- constant(常量~), 248, 300
- linear(线性~), 301

Policy(策略), 451

Polychotomizer(多分器), 53

Polynomial regression(多项式回归), 75

- multivariate(多元~), 104

- Polytree(多树), 407
 POMDP, *see* Partially observable Markov decision process 参见部分可观测的马尔可夫决策过程
 Positive examples(正实例), 21
 Posterior Probability distribution (后验概率分布), 341
 Posterior probability of a class(类后验概率), 50
 Posterior probability of a parameter(参数的后验概率), 67
 Postpruning(后剪枝), 194
 Potential function(势函数), 212, 411
 Power function(功效函数), 498
 Predicate(谓词), 201
 Prediction(预测), 5
 Prepruning(先剪枝), 194
 Principal components analysis(主成分分析), 113
 Prior knowledge(先验知识), 294
 Prior probability distribution(先验概率分布), 341
 Prior probability of a class(类先验概率), 50
 Prior probability of a parameter(参数先验概率), 67
 Probabilistic networks(概率网络), 387
 Probabilistic PCA(概率PCA), 123
 Probably approximately correct learning(概率近似正确学习), 29
 Probit function(概率单位函数) 355
 Product term(乘积项), 211
 Projection pursuit(投影追踪), 274
 Proportion of variance(方差比例), 116
 Propositional rule(命题规则), 201
 Pruning(剪枝)
 postpruning(后~), 194
 prepruning(先~), 194
 set(~集), 194
- ### Q
- Q learning(Q学习), 458
 Quadratic discriminant(二次判别式), 95, 211
 Quantization(量化), 146
- ### R
- Radial basis function(径向基函数), 290
 Random subspace(随机子空间), 421
 Randomization(随机化), 482
 RBF, *see* Radial basis function 参见径向基函数
 Real time recurrent learning(实时递归学习), 272
 Recall(召回), 492
 Receiver operating characteristics(接受者操作曲线), 490
 Receptive field(接收域), 288
 Reconstruction error(重构误差), 119, 146
 Recurrent network(递归网络), 271
 Reference vector(参考向量), 164
 Regression(回归), 9, 35
 linear(线性~), 74
 polynomial(多项式~), 75
 polynomial multivariate(多项式多元~), 104
 robust(鲁棒的~), 329
 Regression tree(回归树), 192
 Regressogram(回归图), 175
 Regularization(正则化), 80, 266
 Regularized discriminant analysis(正则线性判别分析), 100
 Reinforcement learning(增强学习), 13
 Reject(拒绝), 34, 52
 Relative square error(相对平方误差), 76
 Replication(复制), 482
 Representation(表示), 21
 distributed vs. local(分布与局部~), 288
 Ridge regression(岭回归), 266, 350
 Ripper(Ripper算法), 199
 Risk function(风险函数), 51
 Robust regression(鲁棒的回归), 329
 ROC, *see* Receiver operating characteristics 参见接受者操作曲线
 RSE, *see* Relative square error 参见相对平方误差
 Rule(规则)
 extraction(~提取), 295
 induction(~归纳), 198
 pruning(~剪枝), 198
 Rule support(规则的支持度), 198
 Rule value metric(规则价值度量), 199
 Running smoother(移动光滑)
 line(线性~), 177
 mean(均值~), 175

S

Sammon mapping(Sammon 映射), 128
 using MLP(使用 MLP ~), 269
Sammon stress(Sammon 应力), 128
Sample(样本), 48
 correlation(~ 相关性), 89
 covariance(~ 协方差), 89
 mean(~ 均值), 89
Sarsa(Sarsa(算法)), 458
 Sarsa(λ)(Sarsa(λ), (算法)), 461
Scatter(散布), 129
Scree graph(斜坡图), 116
Self-organizing map(自组织映射), 286
Semiparametric density estimation(半参数密度估计), 144
Sensitivity(灵敏性), 493
Sensor fusion(传感器融合), 421
Sequential covering(顺序覆盖), 199
Sigmoid(S 形), 218
Sign test(符号检验), 509
Single-link clustering(单链接聚类), 157
Slack variable(松弛变量), 315
Smoother(光滑子/方法), 174
Smoothing splines(光滑样条), 178
Soft count(软计数), 376
Soft error(软误差), 315
Soft weight sharing(软权重共享), 267
Softmax(软最大), 224
SOM, *see* Self-organizing map 参见自组织映射
Spam filtering(垃圾邮件过滤), 103
Specificity(特效性), 493
Spectral decomposition(谱分解), 115
Speech recognition(语音识别), 380
Sphere node(球形节点), 203
Stability-plasticity dilemma(稳定性-可塑性两难选择), 281
Stacked generalization(层叠泛化), 435
Statlib repository(Statlib 知识库), 17
Stochastic automaton(随机自动机), 364
Stochastic gradient descent(随机梯度下降), 241
Stratification(分层), 487

Strong learner(强学习器), 431
Structural adaptation(结构自适应), 263
Structural risk minimization(结构风险最小化), 82
Subset selection(子集选择), 110
Sum-product algorithm(和-积算法), 412
Supervised learning(监督学习), 9
Support of an association rule(关联规则的支持度), 55
Support vector machine(支持向量机), 313
SVM, *see* Support vector machine 参见支持向量机
Switching HMM(开关 HMM), 400
Synapse(突触), 234
Synaptic weight(突触权重), 237

T

t distribution(t 分布), 495
 t test(t 检验), 498
Tangent prop(正切支撑), 263
TD, *see* Temporal difference 参见时间差分
Template matching(模板匹配), 98
Temporal difference(时间差分), 455
 learning(~ 学习), 458
 TD(0)(TD(0)(算法)), 459
 TD-Gammon(TD-Gammon(算法)), 471
Test set(检验集), 40
Threshold(阈值), 212
 function(~ 函数), 238
Time delay neural network(时间延迟神经网络), 270
Topographical map(地形图), 287
Transition probability(转移概率), 364
Traveling salesman problem(旅行商问题), 306
Triple trade-off(三元权衡), 39
Tukeys test(Tukey 检验), 512
Two-sided confidence interval(双侧置信区间), 494
Two-sided test(双侧检验), 497
Type 2 maximum likelihood procedure(2 型最大似然过程), 360
Type I error(第一类误差), 497
Type II error(第二类误差), 497

U

UCI repository(UCI 知识库), 17
Unbiased estimator(无偏估计), 65

Underfitting(欠拟合), 39, 79
Unfolding in time(按时间展开), 272
Unit normal distribution(单位正态分布), 493
Univariate tree(单变量树), 187
Universal approximation(普适近似), 248
Unobservable variable(不可观测的变量), 48
Unstable algorithm(不稳定算法), 430
Utility function(效用函数), 54
Utility theory(效用理论), 54

V

Validation set(验证集), 40
Value iteration(价值迭代), 453
Value of information(信息价值), 464, 469
Vapnik-Chervonenkis (VC) dimension(VC 维), 27
Variance(方差), 66
Vector quantization(向量量化), 146
 supervised(监督的 ~), 300
Version space(解空间), 24
Vigilance(警戒值), 285

Virtual example(虚拟实例), 262
Viterbi algorithm(Viterbi 算法), 374
Voronoi tessellation(Voronoi 图), 172
Voting(投票表决), 424

W

Weak learner(弱学习器), 431
Weight(权重)
 decay(~ 衰减), 263
 sharing(~ 共享), 260
 sharing soft(共享软 ~), 267
 vector(~ 向量), 212
Winner-take-all(胜者全取), 280
Within-class scatter matrix(类内散布矩阵), 130
Wrappers(包装), 138

Z

z, see Unit normal distribution 参见单位正态分布
z-normalization(z-规范化), 91, 526
Zero-one loss(0-1 损失), 51

推荐阅读



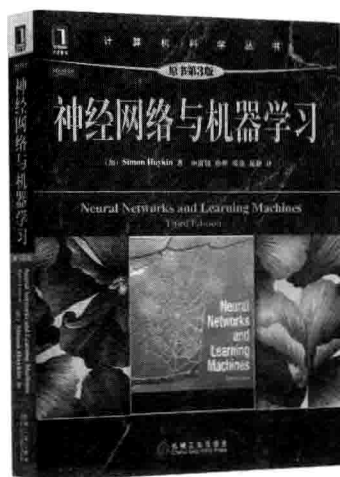
机器学习

作者: (美) Tom Mitchell ISBN: 978-7-111-10993-7 定价: 35.00元



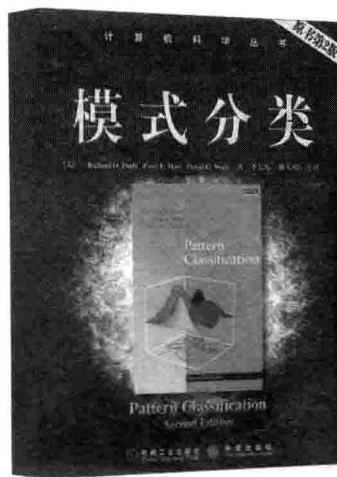
机器学习基础教程

作者: (英) Simon Rogers 等 ISBN: 978-7-111-40702-7 定价: 45.00元



神经网络与机器学习 (原书第3版)

作者: (加) Simon Haykin ISBN: 978-7-111-32413-3 定价: 79.00元



模式分类 (原书第2版)

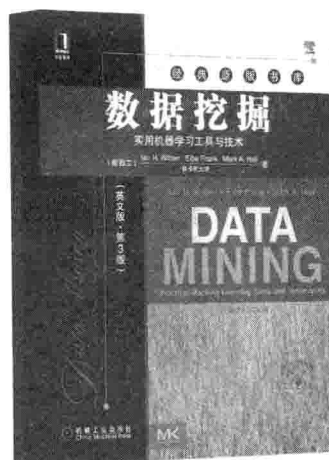
作者: (美) Richard O. Duda 等 ISBN: 978-7-111-12148-1 定价: 59.00元

推荐阅读



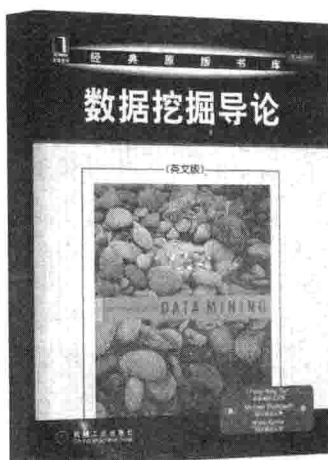
数据挖掘：概念与技术（原书第3版）

作者：（美）Jiawei Han 等
ISBN: 978-7-111-39140-1
定价：79.00元



数据挖掘：实用机器学习工具与技术（英文版·第3版）

作者：（新西兰）Ian H. Witten 等
ISBN: 978-7-111-37417-6
定价：108.00元



数据挖掘导论（英文版）

作者：（美）Pang-Ning Tan 等
ISBN: 978-7-111-31670-1
定价：59.00元



数据挖掘与R语言

作者：（葡）Luis Torgo
ISBN: 978-7-111-40700-3
定价：49.00元

机器学习导论 (原书第2版)

Introduction to Machine Learning Second Edition

机器学习的目标是对计算机编程,以便使用样本数据或以往的经验来解决给定的问题。机器学习已经有许多成功应用,包括分析以往销售数据来预测客户行为,优化机器人的行为以便使用最少的资源来完成任务,以及从生物信息数据中提取知识的各种系统。本书是关于机器学习这一主题内容全面的教科书,涵盖了通常在机器学习导论中并不包括的广泛题材。为了对机器学习问题和解进行统一的论述,本书讨论了源自不同领域的多种方法,包括统计学、模式识别、神经网络、人工智能、信号处理、控制和数据挖掘。书中对所有学习算法都进行了解释,以便读者可以轻易地将书中的公式转变为计算机程序。本书可用作高年级本科生或硕士研究生的教材,也可作为关注机器学习方法应用的专业人员的参考书。

第2版新增内容

- 各章都进行了改写和更新,新增了核机器(第13章)、贝叶斯估计(第14章)和图模型(第16章)内容。
- 在机器学习实验的设计和分析一章(第19章)中扩展了统计检验的内容。
- 在本书配套网站(<http://www.cmpe.boun.edu.tr/~ethem/i2ml2e/>)上提供了PPT、勘误等。
- 增加了一些习题。

作者简介

Ethem Alpaydin 土耳其伊斯坦布尔博阿齐奇大学(Bogaziçi University)计算机工程系教授。他于1990年在瑞士洛桑联邦理工学院获博士学位,1991年在加州大学伯克利分校国际计算机研究所(ICS, UC Berkeley)做博士后工作;之后作为访问学者,先后在美国麻省理工学院、加州大学伯克利分校国际计算机研究所、瑞士戴尔莫尔感知人工智能研究所(IDIAP)从事研究工作。他是土耳其科学院院士,IEEE高级会员,牛津大学出版社《The Computer Journal》杂志编委和Elsevier出版社《Pattern Recognition》杂志副主编。



投稿热线: (010) 88379604
客服热线: (010) 88378991 88361066
购书热线: (010) 68326294 88379649 68995259

华章网站: www.hzbook.com
数字阅读: www.hzmedia.com.cn

上架指导: 计算机/机器学习

ISBN 978-7-111-45377-2



9 787111 453772 >

定价: 59.00元